

# Simple and robust equilibrated flux a posteriori estimates for singularly perturbed reaction–diffusion problems\*

Iain Smears<sup>†</sup>

Martin Vohralík<sup>‡§</sup>

June 2, 2020

## Abstract

We consider energy norm a posteriori error analysis of conforming finite element approximations of singularly perturbed reaction–diffusion problems on simplicial meshes in arbitrary space dimension. Using an equilibrated flux reconstruction, the proposed estimator gives a guaranteed global upper bound on the error without unknown constants, and local efficiency robust with respect to the mesh size and singular perturbation parameters. Whereas previous works on equilibrated flux estimators only considered lowest-order finite element approximations and achieved robustness through the use of boundary-layer adapted submeshes or via combination with residual-based estimators, the present methodology applies in a simple way to arbitrary-order approximations and does not request any submesh or estimators combination. The equilibrated flux is obtained via local reaction–diffusion problems with suitable weights (cut-off factors), and the guaranteed upper bound features the same weights. We prove that the inclusion of these weights is not only sufficient but also necessary for robustness of any flux equilibration estimate that does not employ submeshes or estimators combination, which shows that some of the flux equilibrations proposed in the past cannot be robust. To achieve the fully computable upper bound, we derive explicit bounds for some inverse inequality constants on a simplex, which may be of independent interest.

**Key words:** singular perturbation, a posteriori error analysis, local efficiency, robustness, equilibrated flux

## 1 Introduction

Let  $\Omega$  be a polygonal/polyhedral/polytopal domain in  $\mathbb{R}^d$ ,  $d \geq 1$ , with a Lipschitz-continuous boundary. Let  $\varepsilon > 0$  and  $\kappa \geq 0$  be two fixed real parameters, and let  $f \in L^2(\Omega)$  be a given source term. Consider the problem: find  $u : \Omega \rightarrow \mathbb{R}$  such that

$$-\varepsilon^2 \Delta u + \kappa^2 u = f \quad \text{in } \Omega, \quad (1.1a)$$

$$u = 0 \quad \text{on } \partial\Omega. \quad (1.1b)$$

---

\*This project has received funding from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (grant agreement No 647134 GATIPOR).

<sup>†</sup>Department of Mathematics, University College London, Gower Street, WC1E 6BT London, United Kingdom ([i.smears@ucl.ac.uk](mailto:i.smears@ucl.ac.uk)).

<sup>‡</sup>Inria, 2 rue Simone Iff, 75589 Paris, France ([martin.vohralik@inria.fr](mailto:martin.vohralik@inria.fr)).

<sup>§</sup>Université Paris-Est, CERMICS (ENPC), 77455 Marne-la-Vallée 2, France.

Let  $a(\cdot, \cdot)$  be the symmetric bilinear form defined by

$$a(w, v) := \varepsilon^2(\nabla w, \nabla v) + \kappa^2(w, v), \quad w, v \in H_0^1(\Omega), \quad (1.2)$$

where  $(\cdot, \cdot)$  denotes the  $L^2$ -inner product of scalar- and vector-valued functions on  $\Omega$ , with associated norm  $\|\cdot\|$ . The restriction of the  $L^2$ -inner product to an open subset  $\omega \subset \Omega$  is denoted by  $(\cdot, \cdot)_\omega$ , with associated norm  $\|\cdot\|_\omega$ . The weak formulation of problem (1.1) is to find  $u \in H_0^1(\Omega)$  such that

$$a(u, v) = (f, v) \quad \forall v \in H_0^1(\Omega). \quad (1.3)$$

The energy norm  $\|\cdot\|$  associated to problem (1.1) is then the norm induced by the form  $a(\cdot, \cdot)$ , namely

$$\|v\|^2 := a(v, v), \quad v \in H_0^1(\Omega). \quad (1.4)$$

In this paper, we shall be primarily interested in the case where  $\varepsilon \ll \kappa$ , when problem (1.1) is said to be *singularly perturbed*. Then, the accurate numerical approximation can be challenging due to the typical presence of sharp boundary and/or interior layers in the solution.

In order to present more specifically the focus of this work, let us consider a simplicial mesh  $\mathcal{T}$  of  $\Omega$  and let  $V_{\mathcal{T}} := \mathbb{P}_p(\mathcal{T}) \cap H_0^1(\Omega)$  denote the subspace of  $H_0^1(\Omega)$  of piecewise polynomial functions of degree at most  $p$ , where  $p \geq 1$  is a fixed integer. The conforming Galerkin finite element approximation of (1.3) consists of finding  $u_{\mathcal{T}} \in V_{\mathcal{T}}$  such that

$$a(u_{\mathcal{T}}, v_{\mathcal{T}}) = (f, v_{\mathcal{T}}) \quad \forall v_{\mathcal{T}} \in V_{\mathcal{T}}. \quad (1.5)$$

The goal is to find a computable a posteriori error estimator  $\eta(u_{\mathcal{T}})$  that satisfies

$$\|u - u_{\mathcal{T}}\| \leq C_{\text{rel}}\eta(u_{\mathcal{T}}), \quad \eta(u_{\mathcal{T}}) \leq C_{\text{eff}}\|u - u_{\mathcal{T}}\| + \text{data oscillation}. \quad (1.6)$$

The first inequality in (1.6) is called reliability, while the second inequality is called (global) efficiency. A localized version of the efficiency bound is actually desirable. The quality of the estimator is determined by the product of the two constants  $C_{\text{rel}}$  and  $C_{\text{eff}}$ . A key requirement for singularly perturbed problems is to obtain estimators that are *robust* in the sense that both constants  $C_{\text{rel}}$  and  $C_{\text{eff}}$  are independent of the singular perturbation parameters  $\varepsilon$  and  $\kappa$ . Only such estimates can quantify well the error in the numerical approximation and be reliably used in adaptive algorithms which allow for efficient approximation of the localized features of the solution.

Recently, several methodologies for constructing error estimators that satisfy (1.6) in a robust way have been studied. Verfürth [36] (see also [37] or [39, Section 4.3]) was probably the first to show robust bounds, in the framework of the so-called residual-based estimates. For the problem at hand, these estimators take the form (up to the data oscillation term and possible generic constants)

$$\eta_{\text{res}}(u_{\mathcal{T}})^2 := \sum_{K \in \mathcal{T}} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2 + \sum_{F \in \mathcal{F}_\Omega} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2, \quad (1.7)$$

where the local element and face residuals are defined respectively by

$$r_{\mathcal{T}}|_K := (f + \varepsilon^2 \Delta_{\mathcal{T}} u_{\mathcal{T}} - \kappa^2 u_{\mathcal{T}})|_K, \quad (1.8a)$$

$$j_{\mathcal{T}}|_F := -\varepsilon^2 \llbracket \nabla u_{\mathcal{T}} \cdot \mathbf{n}_F \rrbracket_F, \quad (1.8b)$$

and where  $\Delta_{\mathcal{T}}$  denotes the element-wise Laplacian,  $\llbracket \nabla u_{\mathcal{T}} \cdot \mathbf{n}_F \rrbracket_F$  denotes the jump of the normal component of  $\nabla u_{\mathcal{T}}$  over the face  $F$ ,  $\mathcal{F}_\Omega$  stands for the set of internal faces of the mesh  $\mathcal{T}$ , and the weights (cut-off factors) take the form

$$\alpha_S := \min \left\{ \frac{h_S}{\varepsilon}, \frac{1}{\kappa} \right\}, \quad (1.9)$$

with  $h_S$  being the diameter of  $S$ , where  $S$  is either a simplex  $K$  or a face  $F$ . The resulting estimator  $\eta_{\text{res}}(u_{\mathcal{T}})$  is thus a straightforward extension from the pure diffusion case  $\kappa = 0$  and is simple to implement in practice. The proof that  $\eta_{\text{res}}$  satisfies the second inequality in (1.6) rests on a bubble function technique, where the face bubble functions are defined with respect to a *submesh* matching the boundary-layer length scales and are possibly very steeply decaying. Their role is to capture the sharp layers caused by the singular perturbation. Note that these bubble functions, and hence the submeshes on which they are defined, are only employed in the analysis; thus they do not need to be constructed in practice. Shortly after, Ainsworth and Babuška [2] extended the method of equilibrated residuals, cf. [3], to satisfy (1.6) in a robust way for lowest-order approximations, i.e.  $p = 1$ . In contrast to the residual-based estimators, a boundary-layer adapted submesh in each mesh element needs to be *constructed in practice* in order to evaluate the estimator.

Further progress has been made since, although, to the best of our knowledge, only in the case of lowest-order approximations where the polynomial degree  $p = 1$ . Robust estimates that are guaranteed ( $C_{\text{rel}} = 1$ ) and where  $\eta(u_{\mathcal{T}})$  is fully computable have been obtained in Cheddadi *et al.* [10]. This remedies that  $C_{\text{rel}}$  is unknown for residual-based estimates and that exact solutions of some infinite-dimensional boundary value problems on each element (which cannot be performed exactly in practice) are required in the equilibrated residuals approach. The estimator in [10] is based on an equilibrated flux  $\boldsymbol{\sigma}_{\mathcal{T}}$  belonging to a discrete subspace of  $\mathbf{H}(\text{div})$  that satisfies the equilibration identity  $\nabla \cdot \boldsymbol{\sigma}_{\mathcal{T}} + \kappa^2 u_{\mathcal{T}} = f_{\mathcal{T}}$ , where  $f_{\mathcal{T}}$  is a piecewise polynomial approximation of  $f$ . The estimator is then composed of terms of the form

$$\min \left\{ \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}\|_K, C\varepsilon^{-\frac{1}{2}} \alpha_F^{\frac{1}{2}} \|j_{\mathcal{T}}\|_{\partial K \setminus \partial \Omega} \right\}.$$

Thus it can be seen as a *combination* between an equilibrated flux estimator for diffusion problems and the residual-based estimator of [36] for reaction–diffusion problems. No submesh is needed for the construction of the estimator. Subsequently, Ainsworth and Vejchodský [4, 5] proceed in two stages. First, equilibrated face fluxes are computed as in [2], and then, equilibrated fluxes are obtained by face liftings, so that the final estimate  $\eta(u_{\mathcal{T}})$  is also fully computable and the first inequality in (1.6) is guaranteed with  $C_{\text{rel}} = 1$ . As in [2], though, boundary-layer adapted submeshes appear in the construction of the estimator.

The use of a submesh complicates the construction and implementation of the equilibrated flux estimators of [4, 5]. Moreover, it is likely to be even more involved when moving beyond lowest-order approximations. In this work, by further developing the idea in [10], we show how to obtain *simple*, i.e. avoiding any submesh, yet *robust* equilibrated flux estimators for *arbitrary-order* approximations. The a posteriori error estimates presented in this paper are based on a locally computable flux  $\boldsymbol{\sigma}_{\mathcal{T}}$  and potential approximation  $\phi_{\mathcal{T}}$ , respectively belonging to discrete subspaces of  $\mathbf{H}(\text{div}, \Omega)$  and  $L^2(\Omega)$  of the current mesh  $\mathcal{T}$ , that satisfy the key *equilibration property*

$$\nabla \cdot \boldsymbol{\sigma}_{\mathcal{T}} + \kappa^2 \phi_{\mathcal{T}} = \Pi_{\mathcal{T}} f, \quad (1.10)$$

where  $\Pi_{\mathcal{T}}: L^2(\Omega) \rightarrow \mathbb{P}_p(\mathcal{T})$  denotes the element-wise  $L^2$ -orthogonal projection operator. Note that  $\Pi_{\mathcal{T}}$  can be computed locally and separately for each element. The upper bound on the error then has the simple form

$$\|u - u_{\mathcal{T}}\|^2 \leq \sum_{K \in \mathcal{T}} \left[ w_K \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}\|_K + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K + \tilde{w}_K \|f - \Pi_{\mathcal{T}} f\|_K \right]^2, \quad (1.11)$$

where  $w_K$  is an elementwise computable *weight* (cut-off factor) such that

$$w_K = \min \left\{ 1, C_* \sqrt{\frac{\varepsilon}{\kappa h_K}} \right\}, \quad \tilde{w}_K = \min \left\{ \frac{h_K}{\pi \varepsilon}, \frac{1}{\kappa} \right\},$$

with a fixed computable constant  $C_*$  given by (2.7); see Theorem 3.1 below for further details. The equilibrated flux  $\sigma_{\mathcal{T}}$  and approximate potential  $\phi_{\mathcal{T}}$  in (1.10), (1.11) are obtained by an extension of the patchwise equilibration of [13, 9], see also [8, 20].

Furthermore, we prove robustness and efficiency of the estimator (1.11) by showing that its local contributions are bounded, up to a constant, by the local residual estimators. More precisely, for each  $K \in \mathcal{T}$ , we show that

$$\begin{aligned} w_K^2 \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K^2 + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K^2 &\lesssim \sum_{K' \in \mathfrak{T}_K} \alpha_{K'}^2 \|r_{\mathcal{T}}\|_{K'}^2 + \sum_{F \in \mathfrak{F}_K} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \\ &\lesssim \sum_{K' \in \mathfrak{T}_K} \left[ \|u - u_{\mathcal{T}}\|_{K'}^2 + \alpha_{K'}^2 \|f - \Pi_{\mathcal{T}} f\|_{K'}^2 \right], \end{aligned} \quad (1.12)$$

where  $\mathfrak{T}_K$  and  $\mathfrak{F}_K$  denote the set of elements and faces in a suitable neighbourhood of  $K$  and

$$\|v\|_K^2 := \varepsilon^2 \|\nabla v\|_K^2 + \kappa^2 \|v\|_K^2 \quad K \in \mathcal{T}, \quad (1.13)$$

see Proposition 4.3 and Theorem 4.4 below for full details. Crucially, the constants hidden in  $\lesssim$  in (1.12) are independent of the mesh-sizes  $h_K$  and problem parameters  $\varepsilon$  and  $\kappa$ , depending only on the shape-regularity of  $\mathcal{T}$ , the space dimension  $d$ , and the polynomial degree  $p$ . Hence, just as for residual-based estimates, equilibrated flux estimates have a *straightforward extension* from the pure diffusion case  $\kappa = 0$ , based on including appropriate weights (cut-off factors) and not requiring computations of quantities over any submesh or combination with the residual estimators. In light of these results, we believe that the claims in [38, 39] of a “structural defect” of the robustness of the equilibrated fluxes estimators are not generally valid.

As a side result, we also prove in Proposition 5.1 that the weights  $w_K$  in (1.11) are *necessary* for robustness of *any* equilibrated flux estimate involving the terms  $\|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K$  whenever  $\sigma_{\mathcal{T}}$  is a piecewise polynomial on  $\mathcal{T}$  (and thus its construction does not involve any submesh), regardless of the precise details of the construction of  $\sigma_{\mathcal{T}}$ . This proves that several flux equilibrations proposed in the past cannot be robust with respect to reaction dominance in general (although in many experiments, no loss of robustness may be numerically observed), including those of Repin and Sauter [31], Ainsworth *et al.* [1], Eigel and Samrowski [17], Eigel and Merdon [16], Vejchodský [34, 35], as well as Ainsworth and Vejchodský [6], where the non-robustness is encapsulated in a non-traditional data-oscillation term.

We only treat isotropic meshes. Results for anisotropic meshes can be found in Kunert [29], Grosman [23], Apel *et al.* [7], Zhao and Chen [41], or Kopteva [26, 27]. Also, we are solely interested in the energy norm. Robust estimates in the maximum norm are obtained in Demlow and Kopteva [12] and, on possibly anisotropic meshes, in Kopteva [25] for  $p = 1$  any in Linss [30] for any order  $p \geq 1$  in one space dimension. We refer to Stevenson [33] for robust convergence, and we refer to Faustmann and Melenk [22] and the references therein for balanced norms. Finally, extensions to variable coefficients  $\varepsilon$  and  $\kappa$  can be treated easily as in [5], whereas inhomogeneous Dirichlet and Neumann boundary conditions, mixed parallelepipedal–simplicial meshes, meshes with hanging nodes, and approximations with varying polynomial degree  $p$  can be treated as in Dolejší *et al.* [15]. In particular, the main idea required to treat the case of approximations with varying polynomial degrees is to construct the vertex-patch contributions of the equilibrated flux with a local patch degree greater than or equal to the maximum polynomial degrees of the finite element approximations over the patch [15]. For meshes with hanging nodes, the equilibration can be taken over an extended patch of elements determined by the support of the associated conforming nodal basis function. See also [18, Section 7] concerning the treatment of meshes with arbitrarily many hanging nodes per face.

## 2 Construction of the equilibrated flux

We present in this section the construction of the equilibrated flux  $\boldsymbol{\sigma}_{\mathcal{T}}$  and of the potential approximation  $\phi_{\mathcal{T}}$ .

### 2.1 Notation

Let  $\mathcal{T}$  be a matching simplicial partition of the domain  $\Omega$ , i.e.,  $\bigcup_{K \in \mathcal{T}} K = \overline{\Omega}$ , any element  $K \in \mathcal{T}$  is a closed simplex (interval when  $d = 1$ , triangle when  $d = 2$ , tetrahedron when  $d = 3$ ), and the intersection of two different simplices is either empty, or a vertex, or their common  $l$ -dimensional face,  $1 \leq l \leq d - 1$ . We denote by  $\vartheta_{\mathcal{T}} > 0$  the shape-regularity parameter of the mesh  $\mathcal{T}$ , i.e.

$$\vartheta_{\mathcal{T}} := \max_{K \in \mathcal{T}} \frac{h_K}{\rho_K}, \quad (2.1)$$

where  $\rho_K$  and  $h_K$  are respectively the diameter of the largest ball contained in  $K$  and of  $K$ . For each element  $K \in \mathcal{T}$  and for a fixed integer  $p \geq 1$ , let  $\mathbb{P}_p(K)$  denote the space of polynomials of total degree at most  $p$  on  $K$ . Let

$$\mathbb{P}_p(\mathcal{T}) := \{v \in L^2(\Omega), v|_K \in \mathbb{P}_p(K) \quad \forall K \in \mathcal{T}\}$$

denote the space of scalar piecewise polynomials of degree at most  $p$  over  $\mathcal{T}$ . Let  $\Pi_{\mathcal{T}}: L^2(\Omega) \rightarrow \mathbb{P}_p(\mathcal{T})$  denote the  $L^2$ -orthogonal projection operator from  $L^2(\Omega)$  onto  $\mathbb{P}_p(\mathcal{T})$ . We additionally consider  $\mathbf{L}^2(\Omega) := L^2(\Omega; \mathbb{R}^d)$  and  $\mathbf{RTN}_p(\mathcal{T}) \subset \mathbf{L}^2(\Omega)$  the piecewise Raviart–Thomas–Nédélec space defined by

$$\begin{aligned} \mathbf{RTN}_p(\mathcal{T}) &:= \{\mathbf{v}_{\mathcal{T}} \in \mathbf{L}^2(\Omega), \mathbf{v}_{\mathcal{T}}|_K \in \mathbf{RTN}_p(K) \quad \forall K \in \mathcal{T}\}, \\ \mathbf{RTN}_p(K) &:= \mathbb{P}_p(K; \mathbb{R}^d) + \mathbb{P}_p(K)\mathbf{x}. \end{aligned} \quad (2.2)$$

For any subset  $S$  of  $\overline{\Omega}$ , let  $h_S$  denote the diameter of  $S$ . Thus, in particular,  $h_K$  denotes the diameter of the element  $K \in \mathcal{T}$ . Let  $\mathcal{V}$  denote the set of vertices of the mesh  $\mathcal{T}$ . It is partitioned into the set of interior vertices  $\mathcal{V}^{\text{int}} := \{\mathbf{a} \in \mathcal{V}, \mathbf{a} \in \Omega\}$ , and boundary vertices  $\mathcal{V}^{\text{ext}} := \mathcal{V} \setminus \mathcal{V}^{\text{int}}$ . For each vertex  $\mathbf{a} \in \mathcal{V}$ , the function  $\psi_{\mathbf{a}}$  is the hat function associated with  $\mathbf{a}$ , i.e.,  $\psi_{\mathbf{a}} \in \mathbb{P}_1(\mathcal{T}) \cap H^1(\Omega)$  taking value 1 in the vertex  $\mathbf{a}$  and 0 in the other vertices. The set  $\omega_{\mathbf{a}}$  is the interior of the support of  $\psi_{\mathbf{a}}$  with associated diameter  $h_{\omega_{\mathbf{a}}}$ . Furthermore, let  $\mathcal{T}_{\mathbf{a}}$  denote the restriction of the mesh  $\mathcal{T}$  to  $\omega_{\mathbf{a}}$ , and let  $\mathcal{F}_{\mathbf{a}}$  denote the set of interior faces of  $\mathcal{T}_{\mathbf{a}}$ , i.e. the faces of  $\mathcal{T}_{\mathbf{a}}$  that contain the vertex  $\mathbf{a}$  for  $\mathbf{a} \in \mathcal{V}^{\text{int}}$ , without those on  $\partial\Omega$  for  $\mathbf{a} \in \mathcal{V}^{\text{ext}}$ . For each element  $K \in \mathcal{T}$ , we collect in  $\mathcal{V}_K$  the set of vertices of  $\mathcal{V}$  belonging to  $K$ . We also define  $\mathfrak{T}_K := \bigcup_{\mathbf{a} \in \mathcal{V}_K} \mathcal{T}_{\mathbf{a}}$  and  $\mathfrak{F}_K := \bigcup_{\mathbf{a} \in \mathcal{V}_K} \mathcal{F}_{\mathbf{a}}$ .

Throughout this work, the notation  $a \lesssim b$  means that  $a \leq Cb$  with a constant  $C$  that only depends on the shape-regularity parameter  $\vartheta_{\mathcal{T}}$  of  $\mathcal{T}$ , on the space dimension  $d$ , and on the polynomial degree  $p$ , so that it is in particular independent of the mesh-sizes  $h_K$  and of the problem parameters  $\varepsilon$  and  $\kappa$ ;  $a \simeq b$  then stands for  $a \lesssim b$  and simultaneously  $b \lesssim a$ .

### 2.2 Trace and inverse inequalities

We first recall two inequalities that we will rely on.

**Lemma 2.1** (Trace inequality with explicit constant). *For all  $K \in \mathcal{T}$  and for all  $v \in H^1(K)$  that satisfy  $(v, 1)_K = 0$ , i.e., that have vanishing mean-value on  $K$ , there holds*

$$\|v\|_{\partial K} \leq C_{\text{Tr}} \|\nabla v\|_K^{\frac{1}{2}} \|v\|_K^{\frac{1}{2}}, \quad C_{\text{Tr}} := \sqrt{\vartheta_{\mathcal{T}}(d+1)(2+d/\pi)}. \quad (2.3)$$

*Proof.* We refer the reader to [14, Lemma 1.49] for the explicit constants of the trace inequality for general functions in  $H^1(K)$ ; namely, for each face  $F \subset \partial K$ ,

$$\|v\|_F^2 \leq \vartheta_{\mathcal{T}} (2\|\nabla v\|_K + d/h_K \|v\|_K) \|v\|_K.$$

Then, we additionally apply the Poincaré inequality  $\|v\|_K \leq h_K/\pi \|\nabla v\|_K$  for functions with vanishing mean-value on  $K$ , and sum over all the faces  $F$  to obtain (2.3).  $\square$

**Lemma 2.2** (Inverse inequalities with explicit constants). *For any  $K \in \mathcal{T}$  and any  $\mathbf{v} \in \mathbf{RTN}_p(K)$ , we have*

$$h_K^{1/2} \|\mathbf{v} \cdot \mathbf{n}\|_{\partial K} \leq C_{\text{inv},p,\partial} \|\mathbf{v}\|_K, \quad h_K \|\nabla \cdot \mathbf{v}\|_K \leq C_{\text{inv},p} \|\mathbf{v}\|_K, \quad (2.4)$$

where the constants  $C_{\text{inv},p,\partial}$  and  $C_{\text{inv},p}$  are given by

$$C_{\text{inv},p,\partial} := \sqrt{(d+1)(p+2)(p+d+1)} \vartheta_{\mathcal{T}}, \quad (2.5)$$

$$C_{\text{inv},p} := \sqrt{d} \vartheta_{\mathcal{T}} \frac{\sqrt{5}}{4} (2\sqrt{2})^d \sqrt{(p+1)(p+2)(p+3)(p+4)}. \quad (2.6)$$

*Proof.* See Appendix A.  $\square$

In practice, possibly sharper constants can be obtained for the inequalities in (2.4) by solving numerically small eigenvalue problems on each mesh element, or on a reference element in combination with bounds for the influence of the affine mapping.

We will need below the following constant composed of the constants of the trace and inverse inequalities (2.3) and (2.4):

$$C_* := \frac{1}{\sqrt{2}} \left( \frac{1}{\sqrt{\pi}} C_{\text{inv},p} + C_{\text{Tr}} C_{\text{inv},p,\partial} \right). \quad (2.7)$$

### 2.3 Equilibrated flux $\sigma_{\mathcal{T}}$ and postprocessed potential $\phi_{\mathcal{T}}$

The construction of the auxiliary variables  $\sigma_{\mathcal{T}}$  and  $\phi_{\mathcal{T}}$  giving the equilibration (1.10) is based on independent local mixed finite element approximations of residual problems over the patches of elements around mesh vertices.

For each  $\mathbf{a} \in \mathcal{V}$ , let  $\mathbb{P}_p(\mathcal{T}_{\mathbf{a}})$ , respectively  $\mathbf{RTN}_p(\mathcal{T}_{\mathbf{a}})$ , be the restriction of the space  $\mathbb{P}_p(\mathcal{T})$ , respectively  $\mathbf{RTN}_p(\mathcal{T})$ , to the patch of elements  $\mathcal{T}_{\mathbf{a}}$  around the vertex  $\mathbf{a}$ . The local mixed finite element spaces  $\mathbf{V}_{\mathcal{T}}^{\mathbf{a}}$  and  $Q_{\mathcal{T}}^{\mathbf{a}}$  are defined by

$$\mathbf{V}_{\mathcal{T}}^{\mathbf{a}} := \begin{cases} \{\mathbf{v}_{\mathcal{T}} \in \mathbf{H}(\text{div}, \omega_{\mathbf{a}}) \cap \mathbf{RTN}_p(\mathcal{T}_{\mathbf{a}}), \mathbf{v}_{\mathcal{T}} \cdot \mathbf{n} = 0 \text{ on } \partial\omega_{\mathbf{a}}\} & \text{if } \mathbf{a} \in \mathcal{V}^{\text{int}}, \\ \{\mathbf{v}_{\mathcal{T}} \in \mathbf{H}(\text{div}, \omega_{\mathbf{a}}) \cap \mathbf{RTN}_p(\mathcal{T}_{\mathbf{a}}), \mathbf{v}_{\mathcal{T}} \cdot \mathbf{n} = 0 \text{ on } \partial\omega_{\mathbf{a}} \setminus \partial\Omega\} & \text{if } \mathbf{a} \in \mathcal{V}^{\text{ext}}, \end{cases} \quad (2.8a)$$

$$Q_{\mathcal{T}}^{\mathbf{a}} := \begin{cases} \mathbb{P}_p(\mathcal{T}_{\mathbf{a}}) & \text{if } \kappa > 0 \text{ or } \mathbf{a} \in \mathcal{V}^{\text{ext}}, \\ \{q_{\mathcal{T}} \in \mathbb{P}_p(\mathcal{T}_{\mathbf{a}}), (q_{\mathcal{T}}, 1)_{\omega_{\mathbf{a}}} = 0\} & \text{if } \kappa = 0 \text{ and } \mathbf{a} \in \mathcal{V}^{\text{int}}, \end{cases} \quad (2.8b)$$

see Figure 1.

Recall that  $u_{\mathcal{T}} \in V_{\mathcal{T}}$  with  $V_{\mathcal{T}} = \mathbb{P}_p(\mathcal{T}) \cap H_0^1(\Omega)$  is the finite element solution given by (1.5). Let  $C_*$  be the constant composed of the constants of the trace and inverse inequalities and given by (2.7). Our construction is:

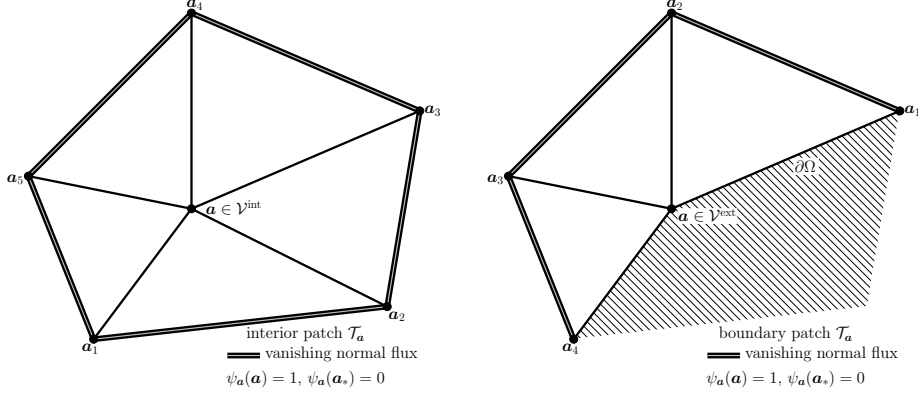


Figure 1: Patches of elements  $\mathcal{T}_a$ , vanishing normal flux conditions in the local Raviart–Thomas–Nédélec spaces  $V_{\mathcal{T}}^a$ , and hat functions  $\psi_a$ : interior (left) and boundary (right) vertex  $a \in \mathcal{V}$

**Definition 2.3** (Flux  $\sigma_{\mathcal{T}}$  and potential  $\phi_{\mathcal{T}}$ ). *For each vertex  $a \in \mathcal{V}$ , let  $(\sigma_{\mathcal{T}}^a, \phi_{\mathcal{T}}^a) \in V_{\mathcal{T}}^a \times Q_{\mathcal{T}}^a$  be defined by the local constrained minimization problem*

$$(\sigma_{\mathcal{T}}^a, \phi_{\mathcal{T}}^a) := \arg \min_{\substack{(v_{\mathcal{T}}, q_{\mathcal{T}}) \in V_{\mathcal{T}}^a \times Q_{\mathcal{T}}^a \\ \nabla \cdot v_{\mathcal{T}} + \kappa^2 q_{\mathcal{T}} = \Pi_{\mathcal{T}}(f\psi_a) - \varepsilon^2 \nabla u_{\mathcal{T}} \cdot \nabla \psi_a}} w_a^2 \|\varepsilon \psi_a \nabla u_{\mathcal{T}} + \varepsilon^{-1} v_{\mathcal{T}}\|_{\omega_a}^2 + \|\kappa [\Pi_{\mathcal{T}}(\psi_a u_{\mathcal{T}}) - q_{\mathcal{T}}]\|_{\omega_a}^2 \quad (2.9a)$$

with the weight

$$w_a := \min \left\{ 1, C_* \sqrt{\frac{\varepsilon}{\kappa h_{\omega_a}}} \right\}. \quad (2.9b)$$

Then, extending each  $\sigma_{\mathcal{T}}^a$  and  $\phi_{\mathcal{T}}^a$  by zero outside of the patch subdomain  $\omega_a$ ,  $\sigma_{\mathcal{T}} \in \mathbf{RTN}_p(\mathcal{T})$  and  $\phi_{\mathcal{T}} \in \mathbb{P}_p(\mathcal{T})$  are given by

$$\sigma_{\mathcal{T}} := \sum_{a \in \mathcal{V}} \sigma_{\mathcal{T}}^a, \quad \phi_{\mathcal{T}} := \sum_{a \in \mathcal{V}} \phi_{\mathcal{T}}^a. \quad (2.9c)$$

We remark that for an interior vertex  $a \in \mathcal{V}^{\text{int}}$ , we have

$$(\Pi_{\mathcal{T}}(f\psi_a) - \varepsilon^2 \nabla u_{\mathcal{T}} \cdot \nabla \psi_a, 1)_{\omega_a} = (f, \psi_a)_{\omega_a} - \varepsilon^2 (\nabla u_{\mathcal{T}}, \nabla \psi_a)_{\omega_a} = \kappa^2 (u_{\mathcal{T}}, \psi_a)_{\omega_a} \quad (2.10)$$

by Galerkin orthogonality with  $\psi_a \in V_{\mathcal{T}}$  as a test function in (1.5). Since

$$(\nabla \cdot \sigma_{\mathcal{T}}^a, 1)_{\omega_a} = (\sigma_{\mathcal{T}}^a \cdot \mathbf{n}_{\omega_a}, 1)_{\partial \omega_a} = 0$$

by Green's theorem and the vanishing normal flux condition imposed in the definition (2.8a) of  $V_{\mathcal{T}}^a$ , it follows that  $\phi_{\mathcal{T}}^a$  necessarily satisfies the mean-value property

$$(\phi_{\mathcal{T}}^a, 1)_{\omega_a} = (\psi_a u_{\mathcal{T}}, 1)_{\omega_a} \quad \forall a \in \mathcal{V}^{\text{int}}$$

whenever  $\kappa > 0$ . If  $\kappa = 0$  instead, then  $\phi_{\mathcal{T}}^a$  is undefined by (2.9a) but one remarks that it is no longer needed anywhere in the paper. In this case, Definition 2.3 coincides with [8, equation (9)], [19, Definition 6.9], or [20, Construction 3.4]; in particular, the Neumann compatibility condition of problem (2.9a) for  $a \in \mathcal{V}^{\text{int}}$  follows from (2.10).

In practice, the constrained minimization problem (2.9a) is solved through its Euler–Lagrange equations, which can be reduced to solving a linear system of dimension  $\dim \mathbf{V}_\mathcal{T}^\mathbf{a} + \dim Q_\mathcal{T}^\mathbf{a}$  in the present context. This problem reads: find  $(\boldsymbol{\sigma}_\mathcal{T}^\mathbf{a}, \phi_\mathcal{T}^\mathbf{a}) \in \mathbf{V}_\mathcal{T}^\mathbf{a} \times Q_\mathcal{T}^\mathbf{a}$  with  $\phi_\mathcal{T}^\mathbf{a} = \gamma_\mathcal{T}^\mathbf{a} + \Pi_\mathcal{T}(\psi_\mathbf{a} u_\mathcal{T})$  and  $(\boldsymbol{\sigma}_\mathcal{T}^\mathbf{a}, \gamma_\mathcal{T}^\mathbf{a}) \in \mathbf{V}_\mathcal{T}^\mathbf{a} \times Q_\mathcal{T}^\mathbf{a}$  such that

$$\varepsilon^{-2} w_\mathbf{a}^2 (\boldsymbol{\sigma}_\mathcal{T}^\mathbf{a}, \mathbf{v}_\mathcal{T})_{\omega_\mathbf{a}} - (\gamma_\mathcal{T}^\mathbf{a}, \nabla \cdot \mathbf{v}_\mathcal{T})_{\omega_\mathbf{a}} = -w_\mathbf{a}^2 (\psi_\mathbf{a} \nabla u_\mathcal{T}, \mathbf{v}_\mathcal{T})_{\omega_\mathbf{a}} \quad \forall \mathbf{v}_\mathcal{T} \in \mathbf{V}_\mathcal{T}^\mathbf{a}, \quad (2.11a)$$

$$(\nabla \cdot \boldsymbol{\sigma}_\mathcal{T}^\mathbf{a}, q_\mathcal{T})_{\omega_\mathbf{a}} + \kappa^2 (\gamma_\mathcal{T}^\mathbf{a}, q_\mathcal{T})_{\omega_\mathbf{a}} = (f \psi_\mathbf{a} - \kappa^2 \psi_\mathbf{a} u_\mathcal{T} - \varepsilon^2 \nabla u_\mathcal{T} \cdot \nabla \psi_\mathbf{a}, q_\mathcal{T})_{\omega_\mathbf{a}} \quad \forall q_\mathcal{T} \in Q_\mathcal{T}^\mathbf{a}. \quad (2.11b)$$

## 2.4 Properties of $\boldsymbol{\sigma}_\mathcal{T}$ and $\phi_\mathcal{T}$

We have constructed  $\boldsymbol{\sigma}_\mathcal{T}$  and  $\phi_\mathcal{T}$  such that the following holds:

**Proposition 2.4** ( *$\mathbf{H}(\text{div}, \Omega)$ -conformity of  $\boldsymbol{\sigma}_\mathcal{T}$ , equilibration*). *Let  $\boldsymbol{\sigma}_\mathcal{T} \in \mathbf{RTN}_p(\mathcal{T})$  and  $\phi_\mathcal{T} \in \mathbb{P}_p(\mathcal{T})$  be given by Definition 2.3. Then  $\boldsymbol{\sigma}_\mathcal{T}$  belongs to  $\mathbf{H}(\text{div}, \Omega)$ , and  $\boldsymbol{\sigma}_\mathcal{T}$  and  $\phi_\mathcal{T}$  satisfy the equilibration property (1.10).*

*Proof.* First, the  $\mathbf{H}(\text{div}, \Omega)$ -conformity of  $\boldsymbol{\sigma}_\mathcal{T}$  follows from the fact that, for any vertex  $\mathbf{a} \in \mathcal{V}$ , the zero extension of  $\boldsymbol{\sigma}_\mathcal{T}^\mathbf{a}$  belongs to  $\mathbf{H}(\text{div}, \Omega)$  as a result of the vanishing normal flux boundary conditions in the space  $\mathbf{V}_\mathcal{T}^\mathbf{a}$ . Then, to show (1.10), we employ the constraint in (2.9a) together with (2.9c):

$$\nabla \cdot \boldsymbol{\sigma}_\mathcal{T} + \kappa^2 \phi_\mathcal{T} = \sum_{\mathbf{a} \in \mathcal{V}} [\nabla \cdot \boldsymbol{\sigma}_\mathcal{T}^\mathbf{a} + \kappa^2 \phi_\mathcal{T}^\mathbf{a}] = \sum_{\mathbf{a} \in \mathcal{V}} [\Pi_\mathcal{T}(f \psi_\mathbf{a}) - \varepsilon^2 \nabla u_\mathcal{T} \cdot \nabla \psi_\mathbf{a}] = \Pi_\mathcal{T} f,$$

where we have used the fact that the hat functions  $\{\psi_\mathbf{a}\}_{\mathbf{a} \in \mathcal{V}}$  form a partition of unity over  $\Omega$ , i.e.  $\sum_{\mathbf{a} \in \mathcal{V}} \psi_\mathbf{a} = 1$ .  $\square$

## 3 A computable guaranteed a posteriori error estimate

This section presents our guaranteed and fully computable a posteriori error estimate. The following upper bound on the energy norm of the error builds on [10, Theorems 3.1 and 4.4] and [5, Lemma 2]. It employs additionally the concept of a potential reconstruction  $\phi_\mathcal{T}$  that will turn out crucial for a simple and robust flux equilibration. Moreover, it relies on the trace and inverse inequalities of Section 2.2 to make appear the crucial weights (cut-off factors), with the constant  $C_*$  given by (2.7).

**Theorem 3.1** (Guaranteed a posteriori error estimate). *Let  $u$  be the weak solution of problem (1.1) given by (1.3) and let  $u_\mathcal{T} \in V_\mathcal{T}$  be its finite element approximation given by (1.5). Let  $\boldsymbol{\sigma}_\mathcal{T} \in \mathbf{RTN}_p(\mathcal{T}) \cap \mathbf{H}(\text{div}, \Omega)$  and  $\phi_\mathcal{T} \in \mathbb{P}_p(\mathcal{T})$  be given by Definition 2.3. Then the following upper bound for the energy norm of the error holds:*

$$\|u - u_\mathcal{T}\|^2 \leq \sum_{K \in \mathcal{T}} [w_K \|\varepsilon \nabla u_\mathcal{T} + \varepsilon^{-1} \boldsymbol{\sigma}_\mathcal{T}\|_K + \|\kappa(u_\mathcal{T} - \phi_\mathcal{T})\|_K + \tilde{w}_K \|f - \Pi_\mathcal{T} f\|_K]^2, \quad (3.1)$$

where the weights  $w_K$  and  $\tilde{w}_K$  are respectively defined by

$$w_K := \min \left\{ 1, C_* \sqrt{\frac{\varepsilon}{\kappa h_K}} \right\}, \quad \tilde{w}_K := \min \left\{ \frac{h_K}{\pi \varepsilon}, \frac{1}{\kappa} \right\}, \quad K \in \mathcal{T}. \quad (3.2)$$



*Proof.* First, we note that the energy norm of the error  $\|u - u_{\mathcal{T}}\|$  is related to the residual  $\mathcal{R}(u_{\mathcal{T}}) \in H^{-1}(\Omega)$ , defined by

$$\langle \mathcal{R}(u_{\mathcal{T}}), v \rangle := (f, v) - a(u_{\mathcal{T}}, v), \quad v \in H_0^1(\Omega),$$

through the identity

$$\|u - u_{\mathcal{T}}\| = \|R(u_{\mathcal{T}})\|_*, \quad \|R(u_{\mathcal{T}})\|_* := \sup_{v \in H_0^1(\Omega), \|v\|=1} \langle \mathcal{R}(u_{\mathcal{T}}), v \rangle, \quad (3.3)$$

cf., e.g., [36, equation (4.1)]. Consider now  $\langle \mathcal{R}(u_{\mathcal{T}}), v \rangle$  for a fixed function  $v \in H_0^1(\Omega)$ . Since  $\sigma_{\mathcal{T}} \in \mathbf{H}(\text{div}, \Omega)$  and  $v \in H_0^1(\Omega)$ , Green's theorem gives  $(\sigma_{\mathcal{T}}, \nabla v) + (\nabla \cdot \sigma_{\mathcal{T}}, v) = 0$ , so

$$\langle \mathcal{R}(u_{\mathcal{T}}), v \rangle = (f, v) - a(u_{\mathcal{T}}, v) = (f - \Pi_{\mathcal{T}} f, v) + (\kappa(\phi_{\mathcal{T}} - u_{\mathcal{T}}), \kappa v) - (\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}, \varepsilon \nabla v), \quad (3.4)$$

where we have also used the equilibration identity (1.10). We now proceed by estimating each term in (3.4) elementwise.

For each element  $K \in \mathcal{T}$ , we use the identity  $(f - \Pi_{\mathcal{T}} f, v)_K = (f - \Pi_{\mathcal{T}} f, v - \Pi_{\mathcal{T}} v)_K$  and the Poincaré–Friedrichs inequality on the convex element  $K$ , i.e.  $\|v - \Pi_{\mathcal{T}} v\|_K \leq \frac{h_K}{\pi} \|\nabla v\|_K$  for any  $v \in H^1(K)$ , together with the energy error definition (1.13), to obtain the following bound

$$|(f - \Pi_{\mathcal{T}} f, v)_K| \leq \|f - \Pi_{\mathcal{T}} f\|_K \min \left\{ \frac{h_K}{\pi \varepsilon} \|\varepsilon \nabla v\|_K, \frac{1}{\kappa} \|\kappa v\|_K \right\} \leq \tilde{w}_K \|f - \Pi_{\mathcal{T}} f\|_K \|v\|_K. \quad (3.5)$$

Here, actually, a little sharper bound is possible by a convex combination of the two possibilities, but we prefer to use the simple form (3.5) with  $\tilde{w}_K$  in the form of minimum given by (3.2).

Next, it is clear that

$$|(\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}, \varepsilon \nabla v)_K| \leq \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K \|v\|_K \quad (3.6)$$

for each  $K \in \mathcal{T}$ . However, this is not necessarily the sharpest possible bound in the singularly perturbed regime  $\kappa \gg \varepsilon$ . Therefore, following the idea of [10, Proof of Theorem 4.4], we use Green's theorem elementwise together with the fact that  $\nabla v = \nabla(v - \bar{v}_K)$ , where  $\bar{v}_K$  denotes the mean-value of  $v$  on  $K$ . This gives

$$(\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}, \varepsilon \nabla v)_K = ((\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}) \cdot \mathbf{n}, \varepsilon(v - \bar{v}_K))_{\partial K} - (\nabla \cdot (\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}), \varepsilon(v - \bar{v}_K))_K.$$

The  $L^2(K)$ -stability of the mean-value,  $\|v - \bar{v}_K\|_K \leq \|v\|_K$ , Young's inequality

$$\|\varepsilon \nabla v\|_K^{\frac{1}{2}} \|\kappa v\|_K^{\frac{1}{2}} \leq \frac{1}{\sqrt{2}} \|v\|_K, \quad (3.7)$$

and the multiplicative trace inequality (2.3) altogether lead to

$$\varepsilon \|v - \bar{v}_K\|_{\partial K} \leq C_{\text{Tr}} \varepsilon^{\frac{1}{2}} \|\varepsilon \nabla v\|_K^{\frac{1}{2}} \|v\|_K^{\frac{1}{2}} \leq \frac{C_{\text{Tr}}}{\sqrt{2}} \sqrt{h_K} \sqrt{\frac{\varepsilon}{\kappa h_K}} \|v\|_K.$$

Combined with the inverse inequality (2.4), we find that

$$|((\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}) \cdot \mathbf{n}, \varepsilon(v - \bar{v}_K))_{\partial K}| \leq \frac{C_{\text{inv}, p, \partial} C_{\text{Tr}}}{\sqrt{2}} \sqrt{\frac{\varepsilon}{\kappa h_K}} \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K \|v\|_K. \quad (3.8)$$

The  $L^2(K)$ -stability of the mean-value, the Poincaré–Friedrichs inequality in the form  $\|v - \bar{v}_K\|_K \leq \frac{h_K}{\pi} \|\nabla v\|_K$ , and (3.7) yield

$$\varepsilon \|v - \bar{v}_K\|_K \leq \varepsilon \|v\|_K^{\frac{1}{2}} h_K^{\frac{1}{2}} \pi^{-\frac{1}{2}} \|\nabla v\|_K^{\frac{1}{2}} \leq \frac{h_K}{\sqrt{2\pi}} \sqrt{\frac{\varepsilon}{\kappa h_K}} \|v\|_K.$$

Thus, combined with the inverse inequality (2.4), we find that

$$|(\nabla \cdot (\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}), \varepsilon(v - \bar{v}_K))_K| \leq \frac{1}{\sqrt{2}} \frac{1}{\sqrt{\pi}} C_{\text{inv},p} \sqrt{\frac{\varepsilon}{\kappa h_K}} \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}\|_K \|v\|_K. \quad (3.9)$$

Therefore, combining inequalities (3.6), (3.8), and (3.9), we get

$$|(\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}, \varepsilon \nabla v)_K| \leq w_K \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}\|_K \|v\|_K \quad \forall K \in \mathcal{T} \quad (3.10)$$

with  $w_K$  given by (3.2) and  $C_*$  given in (2.7). As a side remark, it is possible to obtain a slightly sharper bound, at the expense of making the weight  $w_K$  more complicated than the simple form given by (3.2).

Finally, we can apply the Cauchy–Schwarz inequality to see that  $|(\kappa(\phi_{\mathcal{T}} - u_{\mathcal{T}}), \kappa v)_K| \leq \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K \|v\|_K$ . Therefore, we deduce from (3.4) and the above inequalities that

$$|\langle \mathcal{R}(u_{\mathcal{T}}), v \rangle| \leq \sum_{K \in \mathcal{T}} [w_K \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}\|_K + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K + \tilde{w}_K \|f - \Pi_{\mathcal{T}} f\|_K] \|v\|_K,$$

which implies the upper bound on the error (3.1) after another Cauchy–Schwarz inequality, using (3.3) and  $\sum_{K \in \mathcal{T}} \|v\|_K^2 = \|v\|^2$ .  $\square$

## 4 Efficiency and robustness of the estimate

This section establishes the local (and consequently global) efficiency and robustness of our a posteriori error estimate.

### 4.1 A basic stability result

The main tool in the analysis of efficiency is the following stability result, where, we recall, the broken and the patchwise  $\mathbf{H}(\text{div})$ -conforming Raviart–Thomas–Nédélec spaces  $\mathbf{RTN}_p(\mathcal{T}_{\mathbf{a}})$  and  $\mathbf{V}_{\mathcal{T}}^{\mathbf{a}}$  are respectively given by (2.2) and (2.8).

**Lemma 4.1** (Stability of patchwise flux equilibration). *Let a vertex  $\mathbf{a} \in \mathcal{V}$  be fixed, and let  $g_{\mathcal{T}} \in \mathbb{P}_p(\mathcal{T}_{\mathbf{a}})$  and  $\boldsymbol{\tau}_{\mathcal{T}} \in \mathbf{RTN}_p(\mathcal{T}_{\mathbf{a}})$  be given discontinuous piecewise polynomial functions, with the Neumann compatibility condition  $(g_{\mathcal{T}}, 1)_{\omega_{\mathbf{a}}} = 0$  satisfied if  $\mathbf{a} \in \mathcal{V}^{\text{int}}$ . Then, there holds*

$$\min_{\substack{\mathbf{v}_{\mathcal{T}} \in \mathbf{V}_{\mathcal{T}}^{\mathbf{a}} \\ \nabla \cdot \mathbf{v}_{\mathcal{T}} = g_{\mathcal{T}}}} \|\boldsymbol{\tau}_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}\|_{\omega_{\mathbf{a}}} \lesssim \sup_{\substack{v \in H_*^1(\omega_{\mathbf{a}}) \\ \|\nabla v\|_{\omega_{\mathbf{a}}} = 1}} \{(g_{\mathcal{T}}, v)_{\omega_{\mathbf{a}}} - (\boldsymbol{\tau}_{\mathcal{T}}, \nabla v)_{\omega_{\mathbf{a}}}\}, \quad (4.1)$$

where  $H_*^1(\omega_{\mathbf{a}})$  is the subspace of functions in  $H^1(\omega_{\mathbf{a}})$  that have mean-value zero on the patch subdomain  $\omega_{\mathbf{a}}$  if  $\mathbf{a} \in \mathcal{V}^{\text{int}}$  is an interior vertex, or that vanish on  $\partial\omega_{\mathbf{a}} \cap \partial\Omega$  if  $\mathbf{a} \in \mathcal{V}^{\text{ext}}$  is a boundary vertex.

The above result holds for any dimension  $d \geq 1$ , although some additional properties are known for  $d \leq 3$ . Indeed, in the case where  $d = 2$ , it is shown in [8, Theorem 7] that the constant in (4.1) is in fact independent of the polynomial degree  $p$ , i.e.  $p$ -robust. The extension of the  $p$ -robustness of the bound to the case of  $d = 3$  was shown in [21, Corollaries 3.3 and 3.6]. It is also possible to extend similar results of this kind to situations with hanging nodes and locally refined submeshes, as shown in [18].

## 4.2 Stability with respect to residual estimators

The next lemma shows that the local contributions of the equilibrated flux a posteriori estimators of Definition 2.3 lie below the local residual estimators as defined in (1.7), with the element residuals  $r_{\mathcal{T}}$  and face residuals  $j_{\mathcal{T}}$  are defined by (1.8) and the weights  $\alpha_K$  and  $\alpha_F$  defined by (1.9).

**Lemma 4.2** (Stability of patchwise flux equilibration with respect to residual estimators). *For each  $\mathbf{a} \in \mathcal{V}$ , let  $\boldsymbol{\sigma}_{\mathcal{T}}^{\mathbf{a}}$  and  $\phi_{\mathcal{T}}^{\mathbf{a}}$  be defined by (2.9a). Then*

$$w_{\mathbf{a}}^2 \|\varepsilon \psi_{\mathbf{a}} \nabla u_{\mathcal{T}} + \varepsilon^{-1} \boldsymbol{\sigma}_{\mathcal{T}}^{\mathbf{a}}\|_{\omega_{\mathbf{a}}}^2 + \|\kappa [\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} u_{\mathcal{T}}) - \phi_{\mathcal{T}}^{\mathbf{a}}]\|_{\omega_{\mathbf{a}}}^2 \lesssim \sum_{K \in \mathcal{T}_{\mathbf{a}}} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2 + \sum_{F \in \mathcal{F}_{\mathbf{a}}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2. \quad (4.2)$$

*Proof.* Let a vertex  $\mathbf{a} \in \mathcal{V}$  be fixed. Since  $\boldsymbol{\sigma}_{\mathcal{T}}^{\mathbf{a}}$  and  $\phi_{\mathcal{T}}^{\mathbf{a}}$  are defined as minimizers of the functional in the right-hand side of (2.9a), it is enough to prove that there always exist some  $\mathbf{v}_{\mathcal{T}}^* \in \mathbf{V}_{\mathcal{T}}^{\mathbf{a}}$  and  $q_{\mathcal{T}}^* \in Q_{\mathcal{T}}^{\mathbf{a}}$  that satisfy the constraint  $\nabla \cdot \mathbf{v}_{\mathcal{T}}^* + \kappa^2 q_{\mathcal{T}}^* = \Pi_{\mathcal{T}}(f \psi_{\mathbf{a}}) - \varepsilon^2 \nabla u_{\mathcal{T}} \cdot \nabla \psi_{\mathbf{a}}$  and that satisfy the bound (4.2) with  $\mathbf{v}_{\mathcal{T}}^*$  in place of  $\boldsymbol{\sigma}_{\mathcal{T}}^{\mathbf{a}}$  and  $q_{\mathcal{T}}^*$  in place of  $\phi_{\mathcal{T}}^{\mathbf{a}}$ . The specific construction depends on the mesh size and the problem parameters  $\varepsilon$  and  $\kappa$ , as we now show.

*Case 1,  $\varepsilon/h_{\omega_{\mathbf{a}}} \leq \kappa$  (reaction dominance).* Up to a constant, we have  $\kappa^{-1} \lesssim h_K/\varepsilon$  and  $\kappa^{-1} \lesssim h_F/\varepsilon$  for all elements  $K \in \mathcal{T}_{\mathbf{a}}$  and all interior faces  $F \in \mathcal{F}_{\mathbf{a}}$ . In this case, we adopt the following construction. Let

$$\rho_{\mathbf{a}} := \frac{1}{|\omega_{\mathbf{a}}|} (\psi_{\mathbf{a}} r_{\mathcal{T}}, 1)_{\omega_{\mathbf{a}}} = \frac{1}{|\omega_{\mathbf{a}}|} (\psi_{\mathbf{a}} (f + \varepsilon^2 \Delta_{\mathcal{T}} u_{\mathcal{T}} - \kappa^2 u_{\mathcal{T}}), 1)_{\omega_{\mathbf{a}}}, \quad \mathbf{a} \in \mathcal{V}^{\text{int}},$$

and  $\rho_{\mathbf{a}} := 0$  otherwise. Next, we define

$$q_{\mathcal{T}}^* := \frac{1}{\kappa^2} (\Pi_{\mathcal{T}}(f \psi_{\mathbf{a}}) + \varepsilon^2 \psi_{\mathbf{a}} \Delta_{\mathcal{T}} u_{\mathcal{T}} - \rho_{\mathbf{a}}), \quad \mathbf{v}_{\mathcal{T}}^* := \arg \min_{\substack{\mathbf{v}_{\mathcal{T}} \in \mathbf{V}_{\mathcal{T}}^{\mathbf{a}} \\ \nabla \cdot \mathbf{v}_{\mathcal{T}} = g_{\mathcal{T}}^*}} \|\varepsilon^2 \psi_{\mathbf{a}} \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}\|_{\omega_{\mathbf{a}}}, \quad (4.3)$$

where

$$g_{\mathcal{T}}^* := -\varepsilon^2 (\nabla u_{\mathcal{T}} \cdot \nabla \psi_{\mathbf{a}} + \psi_{\mathbf{a}} \Delta_{\mathcal{T}} u_{\mathcal{T}}) + \rho_{\mathbf{a}}.$$

It is easy to check that if  $\mathbf{a} \in \mathcal{V}^{\text{int}}$ , then  $(g_{\mathcal{T}}^*, 1)_{\omega_{\mathbf{a}}} = 0$ , since the Galerkin orthogonality (take  $v_{\mathcal{T}} = \psi_{\mathbf{a}}$  in (1.5)) implies that

$$(g_{\mathcal{T}}^*, 1)_{\omega_{\mathbf{a}}} = (f, \psi_{\mathbf{a}}) - \varepsilon^2 (\nabla u_{\mathcal{T}}, \nabla \psi_{\mathbf{a}}) - \kappa^2 (u_{\mathcal{T}}, \psi_{\mathbf{a}}) = 0. \quad (4.4)$$

Therefore, it follows that  $q_{\mathcal{T}}^* \in Q_{\mathcal{T}}^{\mathbf{a}}$  and  $\mathbf{v}_{\mathcal{T}}^* \in \mathbf{V}_{\mathcal{T}}^{\mathbf{a}}$  are well-defined and that they satisfy the constraint  $\nabla \cdot \mathbf{v}_{\mathcal{T}}^* + \kappa^2 q_{\mathcal{T}}^* = \Pi_{\mathcal{T}}(f \psi_{\mathbf{a}}) - \varepsilon^2 \nabla u_{\mathcal{T}} \cdot \nabla \psi_{\mathbf{a}}$ .

We now bound  $w_{\mathbf{a}}^2 \|\varepsilon^2 \psi_{\mathbf{a}} \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}^*\|_{\omega_{\mathbf{a}}}^2$  and  $\|\kappa [\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} u_{\mathcal{T}}) - q_{\mathcal{T}}^*]\|_{\omega_{\mathbf{a}}}^2$ . First, we obtain

$$\|\kappa [\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} u_{\mathcal{T}}) - q_{\mathcal{T}}^*]\|_{\omega_{\mathbf{a}}}^2 = \frac{1}{\kappa^2} \|\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} r_{\mathcal{T}}) - \rho_{\mathbf{a}}\|_{\omega_{\mathbf{a}}}^2 \leq \frac{1}{\kappa^2} \|r_{\mathcal{T}}\|_{\omega_{\mathbf{a}}}^2 \lesssim \sum_{K \in \mathcal{T}_{\mathbf{a}}} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2,$$

where we have used the stability of the  $L^2$ -projection (note that  $\rho_{\mathbf{a}}$  is also the mean value of  $\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} r_{\mathcal{T}})$  on  $\omega_{\mathbf{a}}$  for  $\mathbf{a} \in \mathcal{V}^{\text{int}}$ ) and the fact that  $\|\psi_{\mathbf{a}}\|_{\infty, \omega_{\mathbf{a}}} = 1$  to bound  $\|\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} r_{\mathcal{T}}) - \rho_{\mathbf{a}}\|_{\omega_{\mathbf{a}}}$ . Next, we apply Lemma 4.1 to bound  $w_{\mathbf{a}}^2 \|\varepsilon^2 \psi_{\mathbf{a}} \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}^*\|_{\omega_{\mathbf{a}}}^2$ . Note first that for an interior vertex  $\mathbf{a} \in \mathcal{V}^{\text{int}}$ ,  $(\rho_{\mathbf{a}}, v)_{\omega_{\mathbf{a}}} = 0$  for all  $v \in H_*^1(\omega_{\mathbf{a}})$  since  $v \in H_*^1(\omega_{\mathbf{a}})$  implies that  $v$  is orthogonal

to constant functions on  $\omega_a$ . We find that

$$\begin{aligned}
\|\varepsilon^2 \psi_a \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}^*\|_{\omega_a} &\lesssim \sup_{v \in H_*^1(\omega_a), \|\nabla v\|_{\omega_a}=1} \{ (g_{\mathcal{T}}^*, v)_{\omega_a} - (\varepsilon^2 \nabla u_{\mathcal{T}}, \psi_a \nabla v)_{\omega_a} \} \\
&= \sup_{v \in H_*^1(\omega_a), \|\nabla v\|_{\omega_a}=1} \{ -(\varepsilon^2 \nabla u_{\mathcal{T}}, \nabla(\psi_a v))_{\omega_a} - (\varepsilon^2 \Delta_{\mathcal{T}} u_{\mathcal{T}}, \psi_a v)_{\omega_a} \} \quad (4.5) \\
&= \sup_{v \in H_*^1(\omega_a), \|\nabla v\|_{\omega_a}=1} \sum_{F \in \mathcal{F}_a} (j_{\mathcal{T}}, \psi_a v)_F,
\end{aligned}$$

where the last line follows by elementwise integration by parts. It is then straightforward to deduce from the trace inequalities  $\|v\|_F \lesssim h_K^{-\frac{1}{2}} \|v\|_K + \|\nabla v\|_K^{\frac{1}{2}} \|v\|_K^{\frac{1}{2}}$  and the Poincaré–Friedrichs inequality for functions in  $H_*^1(\omega_a)$ , i.e.  $\|v\|_{\omega_a} \lesssim h_{\omega_a} \|\nabla v\|_{\omega_a}$ , that

$$\|\varepsilon^2 \psi_a \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}^*\|_{\omega_a}^2 \lesssim h_{\omega_a} \sum_{F \in \mathcal{F}_a} \|j_{\mathcal{T}}\|_F^2. \quad (4.6)$$

Consequently, using definition (2.9b) of the weight  $w_a$

$$w_a^2 \|\varepsilon \psi_a \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}^*\|_{\omega_a}^2 \lesssim \frac{\varepsilon}{\kappa h_{\omega_a}} \frac{h_{\omega_a}}{\varepsilon^2} \sum_{F \in \mathcal{F}_a} \|j_{\mathcal{T}}\|_F^2 \lesssim \sum_{F \in \mathcal{F}_a} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2.$$

Therefore, if  $\varepsilon/h_{\omega_a} \leq \kappa$ , we have shown that there exist  $\mathbf{v}_{\mathcal{T}}^*$  and  $q_{\mathcal{T}}^*$  satisfying the constraint  $\nabla \cdot \mathbf{v}_{\mathcal{T}}^* + \kappa^2 q_{\mathcal{T}}^* = \Pi_{\mathcal{T}}(f \psi_a) - \varepsilon^2 \nabla u_{\mathcal{T}} \cdot \nabla \psi_a$  and such that

$$w_a^2 \|\varepsilon \psi_a \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}^*\|_{\omega_a}^2 + \|\kappa [\Pi_{\mathcal{T}}(\psi_a u_{\mathcal{T}}) - q_{\mathcal{T}}^*]\|_{\omega_a}^2 \lesssim \sum_{K \in \mathcal{T}_a} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2 + \sum_{F \in \mathcal{F}_a} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2.$$

As explained above, this implies (4.2) in the case  $\varepsilon/h_{\omega_a} \leq \kappa$ .

*Case 2,  $\varepsilon/h_{\omega_a} > \kappa$  (diffusion dominance).* We select

$$q_{\mathcal{T}}^* := \Pi_{\mathcal{T}}(\psi_a u_{\mathcal{T}}), \quad \mathbf{v}_{\mathcal{T}}^* := \arg \min_{\substack{\mathbf{v}_{\mathcal{T}} \in \mathbf{V}_{\mathcal{T}}^a \\ \nabla \cdot \mathbf{v}_{\mathcal{T}} = g_{\mathcal{T}}^*}} \|\varepsilon^2 \psi_a \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}\|_{\omega_a},$$

where

$$g_{\mathcal{T}}^* := \Pi_{\mathcal{T}}(\psi_a (f - \kappa^2 u_{\mathcal{T}})) - \varepsilon^2 \nabla \psi_a \cdot \nabla u_{\mathcal{T}}.$$

Notice that Galerkin orthogonality implies that  $(g_{\mathcal{T}}^*, 1)_{\omega_a} = 0$  if  $\mathbf{a} \in \mathcal{V}^{\text{int}}$  as in (4.4), and also  $\nabla \cdot \mathbf{v}_{\mathcal{T}}^* + \kappa^2 q_{\mathcal{T}}^* = \Pi_{\mathcal{T}}(f \psi_a) - \varepsilon^2 \nabla u_{\mathcal{T}} \cdot \nabla \psi_a$ , so the requested constraint is satisfied. It then follows directly from Lemma 4.1 that

$$\|\varepsilon^2 \psi_a \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}^*\|_{\omega_a} \lesssim \sup_{v \in H_*^1(\omega_a), \|\nabla v\|_{\omega_a}=1} \left\{ (\Pi_{\mathcal{T}}(\psi_a r_{\mathcal{T}}), v)_{\omega_a} + \sum_{F \in \mathcal{F}_a} (\psi_a j_{\mathcal{T}}, v)_F \right\},$$

where we use the fact that elementwise integration by parts shows that, as in (4.5),

$$(g_{\mathcal{T}}^*, v)_{\omega_a} - (\varepsilon^2 \psi_a \nabla u_{\mathcal{T}}, \nabla v)_{\omega_a} = (\Pi_{\mathcal{T}}(\psi_a r_{\mathcal{T}}), v)_{\omega_a} + \sum_{F \in \mathcal{F}_a} (\psi_a j_{\mathcal{T}}, v)_F.$$

Thus, proceeding as in (4.5)–(4.6) for the face residuals term and using the stability of the  $L^2$ -projection,  $\|\psi_a\|_{\infty, \omega_a} = 1$ , and the Poincaré–Friedrichs inequality for functions in  $H_*^1(\omega_a)$ ,  $\|v\|_{\omega_a} \lesssim h_{\omega_a} \|\nabla v\|_{\omega_a}$ , for the element residuals term, we get

$$\|\varepsilon^2 \psi_a \nabla u_{\mathcal{T}} + \mathbf{v}_{\mathcal{T}}^*\|_{\omega_a}^2 \lesssim h_{\omega_a}^2 \sum_{K \in \mathcal{T}_a} \|r_{\mathcal{T}}\|_K^2 + h_{\omega_a} \sum_{F \in \mathcal{F}_a} \|j_{\mathcal{T}}\|_F^2.$$

Consequently,

$$\|\varepsilon\psi_{\mathbf{a}}\nabla u_{\mathcal{T}} + \varepsilon^{-1}\mathbf{v}_{\mathcal{T}}^*\|_{\omega_{\mathbf{a}}}^2 \lesssim \sum_{K \in \mathcal{T}_{\mathbf{a}}} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2 + \sum_{F \in \mathcal{F}_{\mathbf{a}}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2.$$

Hence, on noting that  $w_{\mathbf{a}} \leq 1$  and that  $\|\kappa[\Pi_{\mathcal{T}}(\psi_{\mathbf{a}}u_{\mathcal{T}}) - q_{\mathcal{T}}^*]\|_{\omega_{\mathbf{a}}} = 0$ , we see that (4.2) also holds for the case  $\varepsilon/h_{\omega_{\mathbf{a}}} > \kappa$ .  $\square$

Recall that  $\mathfrak{T}_K := \bigcup_{\mathbf{a} \in \mathcal{V}_K} \mathcal{T}_{\mathbf{a}}$  and  $\mathfrak{F}_K := \bigcup_{\mathbf{a} \in \mathcal{V}_K} \mathcal{F}_{\mathbf{a}}$ .

**Proposition 4.3** (Bound on flux estimators by the residual estimators). *Let  $\sigma_{\mathcal{T}}$  and  $\phi_{\mathcal{T}}$  be given by Definition 2.3. Additionally, let the volume and face residual functions  $r_{\mathcal{T}}$  and  $j_{\mathcal{T}}$  be defined by (1.8). Then, for each element  $K \in \mathcal{T}$ , we have the bound*

$$w_K^2 \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K^2 + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K^2 \lesssim \sum_{K' \in \mathfrak{T}_K} \alpha_{K'}^2 \|r_{\mathcal{T}}\|_{K'}^2 + \sum_{F \in \mathfrak{F}_K} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2. \quad (4.7)$$

*Proof.* For each mesh element  $K \in \mathcal{T}$ , we have  $\sigma_{\mathcal{T}}|_K = \sum_{\mathbf{a} \in \mathcal{V}_K} \sigma_{\mathcal{T}}^{\mathbf{a}}|_K$  and  $\phi_{\mathcal{T}}|_K = \sum_{\mathbf{a} \in \mathcal{V}_K} \phi_{\mathcal{T}}^{\mathbf{a}}|_K$ . Moreover, since  $\{\psi_{\mathbf{a}}\}_{\mathbf{a} \in \mathcal{V}_K}$  form a partition of unity over  $K$  and since  $\Pi_{\mathcal{T}}$  is the elementwise  $L^2$  projection of degree  $p$ , it follows that  $u_{\mathcal{T}}|_K = \Pi_{\mathcal{T}} u_{\mathcal{T}}|_K = \sum_{\mathbf{a} \in \mathcal{V}_K} \Pi_{\mathcal{T}}(\psi_{\mathbf{a}} u_{\mathcal{T}})|_K$ . Furthermore, (3.2) and (2.9b) together with the mesh shape regularity imply that  $w_K \lesssim w_{\mathbf{a}}$  for each  $\mathbf{a} \in \mathcal{V}_K$ , where the constant depends only on  $\vartheta_{\mathcal{T}}$ . Therefore, we obtain

$$\begin{aligned} & w_K^2 \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K^2 + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K^2 \\ & \lesssim \sum_{\mathbf{a} \in \mathcal{V}_K} \left[ w_{\mathbf{a}}^2 \|\varepsilon \psi_{\mathbf{a}} \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}^{\mathbf{a}}\|_K^2 + \|\kappa[\Pi_{\mathcal{T}}(\psi_{\mathbf{a}} u_{\mathcal{T}}) - \phi_{\mathcal{T}}^{\mathbf{a}}]\|_K^2 \right]. \end{aligned}$$

Therefore, we can use (4.2) for each  $\mathbf{a} \in \mathcal{V}_K$  to get (4.7).  $\square$

### 4.3 Local efficiency and robustness of the estimate

We now recall the well-known efficiency and robustness results for residual estimators, see [36, Proposition 4.1] or [39] for details. For each  $K \in \mathcal{T}$  and  $F \in \mathcal{F}_{\Omega}$ , there holds

$$\alpha_K^2 \|r_{\mathcal{T}}\|_K^2 \lesssim \|u - u_{\mathcal{T}}\|_K^2 + \alpha_K^2 \|f - \Pi_{\mathcal{T}} f\|_K^2, \quad (4.8a)$$

$$\varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \lesssim \sum_{K \in \mathcal{T}, F \subset \partial K} \left[ \|u - u_{\mathcal{T}}\|_K^2 + \alpha_K^2 \|f - \Pi_{\mathcal{T}} f\|_K^2 \right]. \quad (4.8b)$$

Therefore, the combination of Proposition 4.3 with (4.8) shows that the equilibrated flux estimator of Theorem 3.1 is locally efficient and robust.

**Theorem 4.4** (Local efficiency and robustness). *Let  $u$  be the weak solution of problem (1.1) given by (1.3) and let  $u_{\mathcal{T}} \in V_{\mathcal{T}}$  be its finite element approximation given by (1.5). Let  $\sigma_{\mathcal{T}} \in \mathbf{RTN}_p(\mathcal{T}) \cap \mathbf{H}(\text{div}, \Omega)$  and  $\phi_{\mathcal{T}} \in \mathbb{P}_p(\mathcal{T})$  be given by Definition 2.3. Then, for each mesh element  $K \in \mathcal{T}$ , there holds*

$$w_K^2 \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K^2 + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K^2 \lesssim \sum_{K' \in \mathfrak{T}_K} \left[ \|u - u_{\mathcal{T}}\|_{K'}^2 + \alpha_{K'}^2 \|f - \Pi_{\mathcal{T}} f\|_{K'}^2 \right], \quad (4.9)$$

where the constant in  $\lesssim$  depends only on the dimension  $d$ , the shape-regularity constant  $\vartheta_{\mathcal{T}}$  of  $\mathcal{T}$ , and on the polynomial degree  $p$ , so that it is independent of the parameters  $\varepsilon$  and  $\kappa$  and the mesh-sizes  $h_K$ .

## 5 Necessity of the weights $w_K$ in the upper bound

Theorems 3.1 and 4.4 show that the estimator  $w_K \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K$  obtained from the flux equilibration of Definition 2.3 is a reliable, locally efficient, and robust energy error estimator for singularly perturbed reaction–diffusion problems. Here we show the *necessity* of the weight  $w_K$  for robustness of equilibrated flux estimators that involve only piecewise polynomial vector fields on  $\mathcal{T}$ . We also recall that an alternative option, related to the approach in [2, 4, 5, 27], is to perform an equilibrations on a submesh.

### 5.1 Necessity of the weights $w_K$

The following proposition applies to any flux equilibration on  $\mathcal{T}$ :

**Proposition 5.1** (Best-possible bound by piecewise polynomials of the mesh  $\mathcal{T}$ ). *Let  $u_{\mathcal{T}} \in \mathbb{P}_p(\mathcal{T}) \cap H_0^1(\Omega)$  be an arbitrary piecewise  $p$ -degree polynomial,  $p \geq 1$ , and let the face residual term  $j_{\mathcal{T}}$  be defined by (1.8b). Let  $\mathbb{P}_{p'}(\mathcal{T}; \mathbb{R}^d)$  denote the space of  $\mathbb{R}^d$ -valued piecewise polynomials of degree at most  $p'$  over  $\mathcal{T}$ , where  $p' \geq 0$  is an arbitrary nonnegative integer. Then,*

$$\inf_{\mathbf{v}_{\mathcal{T}} \in \mathbf{H}(\text{div}, \Omega) \cap \mathbb{P}_{p'}(\mathcal{T}; \mathbb{R}^d)} \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}\| \gtrsim \sqrt{\frac{\kappa \underline{h}}{\varepsilon}} \left( \sum_{F \in \mathcal{F}_{\Omega}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \right)^{\frac{1}{2}}, \quad (5.1)$$

where  $\underline{h} := \min_{K \in \mathcal{T}} h_K$ , and where the constant depends only on the polynomial degrees  $p$  and  $p'$ , the dimension  $d$ , and the shape-regularity  $\vartheta_{\mathcal{T}}$  of  $\mathcal{T}$ .

*Proof.* Let  $\mathbf{v}_{\mathcal{T}} \in \mathbf{H}(\text{div}, \Omega) \cap \mathbb{P}_{p'}(\mathcal{T}; \mathbb{R}^d)$  be arbitrary. Then, for each interior face  $F \in \mathcal{F}_{\Omega}$ , the  $\mathbf{H}(\text{div}, \Omega)$ -conformity of  $\mathbf{v}_{\mathcal{T}}$  implies that  $[\![\mathbf{v}_{\mathcal{T}} \cdot \mathbf{n}_F]\!]_F = 0$ , and hence  $j_{\mathcal{T}}|_F = -\varepsilon [(\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}) \cdot \mathbf{n}_F]_F$ . Since  $(\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}})|_K \in \mathbb{P}_{\max(p', p-1)}(K; \mathbb{R}^d)$  for each element  $K \in \mathcal{T}$ , we can apply the triangle inequality and the inverse inequality (analogous to (2.4)) to find that, for any  $F \in \mathcal{F}_{\Omega}$ ,

$$\varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \lesssim \frac{\varepsilon}{\kappa \underline{h}} \sum_{K \in \mathcal{T}, F \subset \partial K} \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}\|_K^2. \quad (5.2)$$

Therefore, we get (5.1) by summing (5.2) over all faces  $F \in \mathcal{F}_{\Omega}$ , and recalling that  $\mathbf{v}_{\mathcal{T}}$  was arbitrary.  $\square$

The upshot of Proposition 5.1 is that for any problem where the jump estimators are sufficiently dominant, i.e. when

$$\|u - u_{\mathcal{T}}\| \simeq \left( \sum_{F \in \mathcal{F}_{\Omega}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \right)^{\frac{1}{2}}, \quad (5.3)$$

then *any* error estimator involving a term of the form  $\|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}\|$  without any weight will necessarily be non-robust when  $\kappa \underline{h} / \varepsilon$  takes large values, since (5.1) and (5.3) then imply

$$\inf_{\mathbf{v}_{\mathcal{T}} \in \mathbf{H}(\text{div}, \Omega) \cap \mathbb{P}_{p'}(\mathcal{T}; \mathbb{R}^d)} \frac{\|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \mathbf{v}_{\mathcal{T}}\|}{\|u - u_{\mathcal{T}}\|} \gtrsim \sqrt{\frac{\kappa \underline{h}}{\varepsilon}}. \quad (5.4)$$

In other words, the effectivity index can become *arbitrarily large* in the singularly-perturbed regime when the weight  $w_K$  is not included. It is then seen that the inclusion of the weight term  $w_K$  in Theorem 3.1 is *necessary* when considering flux equilibrations from vector-valued

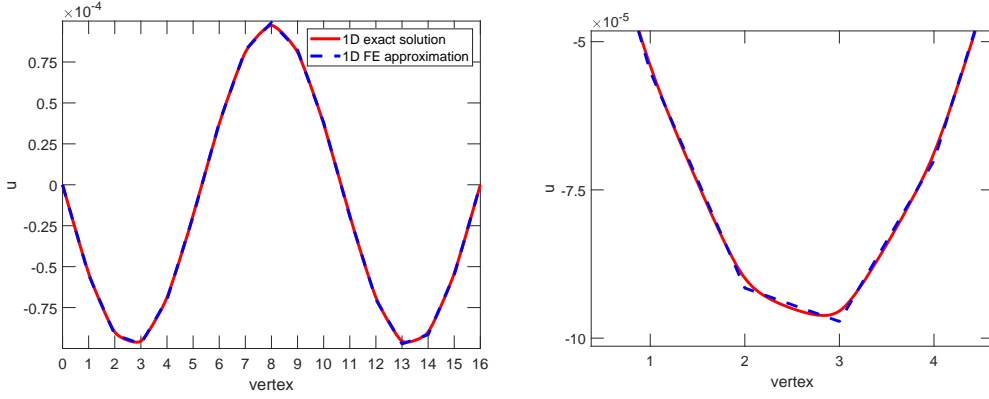


Figure 2: [Example 5.2,  $\varepsilon = 1$ ,  $\kappa = 10^2$ ,  $m = 3$ ] Finite element approximation (5.6) and the exact solution (left), detail (right)

piecewise polynomial subspaces of  $\mathbf{H}(\text{div}, \Omega)$  on the mesh  $\mathcal{T}$ , regardless of the precise details of the construction of the flux. Examples of flux equilibrations proposed in the past that cannot be robust in general include Repin and Sauter [31], Ainsworth *et al.* [1], Eigel and Samrowski [17], Eigel and Merdon [16], Vejchodský [34, 35], as well as Ainsworth and Vejchodský [6], where the non-robustness is encapsulated in a non-traditional data-oscillation term..

We now present an example of a situation where (5.3) holds and where  $\kappa h/\varepsilon$  can be arbitrarily large. In fact the example is similar to the one in [2, Section 2.3], albeit with some suitable adjustments.

**Example 5.2** (Dominant jump estimators). *Let  $\Omega := (-1/2, 1/2)$  and let  $m$  be an odd integer. Consider a uniform mesh  $\mathcal{T}$  of  $\Omega$  with  $2N = (m+1)^2$  intervals,  $N := (m+1)^2/2$ , and mesh size  $h := 1/(2N) = 1/(m+1)^2$ . Hence, the interior nodes are  $x_i = ih$ , where  $i \in \{-N+1, \dots, N-1\}$ . Let*

$$f := \mathcal{I}_{\mathcal{T}} \cos(m\pi x) \in \mathbb{P}_1(\mathcal{T}) \cap H_0^1(\Omega) \quad (5.5)$$

denote the piecewise affine Lagrange interpolant (preserving the point values) of the function  $x \mapsto \cos(m\pi x)$ ; it follows from the fact that  $m$  is odd that  $f \in H_0^1(\Omega)$ . Note that in the example of [2], the function  $f$  was chosen as  $\cos(\pi x)$  instead.

Consider now problem (1.1) along with its finite element approximation (1.5) in the space  $V_{\mathcal{T}} = \mathbb{P}_1(\mathcal{T}) \cap H_0^1(\Omega)$ . It is easy to show that

$$u_{\mathcal{T}} = (\varepsilon^2 \mu_h + \kappa^2)^{-1} f \quad (5.6)$$

is the discrete solution, where

$$\mu_h := \frac{6}{2 + \cos(m\pi h)} \frac{1 - \cos(m\pi h)}{h^2},$$

as a result of the identity

$$\int_{-1/2}^{1/2} f' v'_{\mathcal{T}} dx = \mu_h \int_{-1/2}^{1/2} f v_{\mathcal{T}} dx \quad \forall v_{\mathcal{T}} \in V_{\mathcal{T}}.$$

An illustration of the finite element approximation (5.6) for  $m = 3$  (which gives 16 intervals,  $h = 1/16$ , and  $\mu_h$  roughly equal to  $5.71h^{-1}$ ) together with the exact solution is given in Figure 2.

Now, noting that interior vertices and faces coincide for problems in one space dimension, it is found that

$$r_{\mathcal{T}}|_K = \frac{\varepsilon^2 \mu_h}{\varepsilon^2 \mu_h + \kappa^2} f|_K, \quad j_{\mathcal{T}}|_{x_i} = -\varepsilon^2 \llbracket u'_{\mathcal{T}}(x_i) \rrbracket = \frac{\varepsilon^2}{\varepsilon^2 \mu_h + \kappa^2} \frac{2(1 - \cos(m\pi h))}{h} f(x_i).$$

Moreover, since  $\lim_{m \rightarrow \infty} \frac{1 - \cos(m\pi h)}{h} = \frac{\pi^2}{2}$  when  $h = h(m) = 1/(m+1)^2$ , for any  $m$  there holds  $\mu_h \simeq h^{-1}$ . Suppose also henceforth that  $\kappa h/\varepsilon \geq 1$ , so that  $\alpha_K$  given by (1.9) takes the value  $1/\kappa$ . Then, we find that

$$\sum_{K \in \mathcal{T}} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2 = \frac{1}{\kappa^2} \left( \frac{\varepsilon^2}{\varepsilon^2 \mu_h + \kappa^2} \right)^2 \mu_h^2 \|f\|^2 \simeq \left( \frac{\varepsilon^2}{\varepsilon^2 \mu_h + \kappa^2} \right)^2 \frac{1}{\kappa^2 h^2}.$$

We also obtain

$$\sum_{F \in \mathcal{F}_{\Omega}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \simeq \left( \frac{\varepsilon^2}{\varepsilon^2 \mu_h + \kappa^2} \right)^2 \frac{1}{\varepsilon \kappa} \sum_{i=-N+1}^{N-1} |f(x_i)|^2 \simeq \left( \frac{\varepsilon^2}{\varepsilon^2 \mu_h + \kappa^2} \right)^2 \frac{1}{\varepsilon \kappa h},$$

where we have used the trigonometric identity

$$\sum_{i=-N+1}^{N-1} |f(x_i)|^2 = \sum_{i=-N+1}^{N-1} \left| \cos\left(\frac{m\pi i}{2N}\right) \right|^2 = \sum_{i=-N+1}^{N-1} \left| \cos\left(\frac{(\sqrt{2N}-1)\pi i}{2N}\right) \right|^2 = N = \frac{1}{2h}.$$

Since  $\varepsilon \kappa h \leq \kappa^2 h^2$ , we see that

$$\sum_{K \in \mathcal{T}} \alpha_K^2 \|r_{\mathcal{T}}\|_K^2 \lesssim \sum_{F \in \mathcal{F}_{\Omega}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2 \iff \|u - u_{\mathcal{T}}\|^2 \simeq \sum_{F \in \mathcal{F}_{\Omega}} \varepsilon^{-1} \alpha_F \|j_{\mathcal{T}}\|_F^2,$$

where we note that there is no data oscillation since  $f \in \mathbb{P}_1(\mathcal{T})$ . Hence this provides an example where (5.3) holds, and the factor  $\kappa h/\varepsilon$  can be made arbitrarily large. In Figure 2, the jumps in the derivative of the numerical solution are clearly apparent.

## 5.2 Flux equilibration on a submesh

**Remark 5.3** (Flux equilibration on boundary-layer adapted submeshes). The approach in [4, 5, 27], following [2], can be seen as defining a flux  $\sigma_{\tilde{\mathcal{T}}} \in \mathbf{H}(\text{div}, \Omega)$  that satisfies an equilibration property similar to (1.10), yet with the key difference that  $\sigma_{\tilde{\mathcal{T}}}$  is defined with respect to a submesh  $\tilde{\mathcal{T}}$  of  $\mathcal{T}$  with thin elements that are adapted to the parameters  $\varepsilon$  and  $\kappa$  and local mesh-size (see e.g. [4, Fig. 3]). In this case, the argument in the proof of Proposition 5.1 does not apply, because the inverse inequality  $\|(\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\tilde{\mathcal{T}}}) \cdot \mathbf{n}_F\|_F \lesssim h_K^{-\frac{1}{2}} \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\tilde{\mathcal{T}}}\|_K$ ,  $F \subset \partial K$ ,  $K \in \mathcal{T}$ , is not applicable when  $\sigma_{\tilde{\mathcal{T}}} \in \mathbb{P}_{p'}(\tilde{\mathcal{T}}; \mathbb{R}^d)$  but  $\sigma_{\tilde{\mathcal{T}}} \notin \mathbb{P}_{p'}(\mathcal{T}; \mathbb{R}^d)$ . This essentially shows how there are now two different approaches to constructing robust equilibrated flux estimators. Either the flux is computed as a piecewise polynomial vector field with respect to the original mesh, in which case the inclusion of a weight of the form of  $w_K$  from (3.2) in the upper bound is necessary, or one constructs the flux with respect to some other sufficiently rich subspace of  $\mathbf{H}(\text{div}, \Omega)$ , such as a piecewise polynomial subspace with respect to an adapted submesh  $\tilde{\mathcal{T}}$  of  $\mathcal{T}$ , in which case the weights are not necessary. Note that in some cases, the submeshes are only used conceptually in the derivation of the estimators, with all computations performed in practice without explicitly constructing the submeshes, see, e.g., [27, Lemma 7.1].



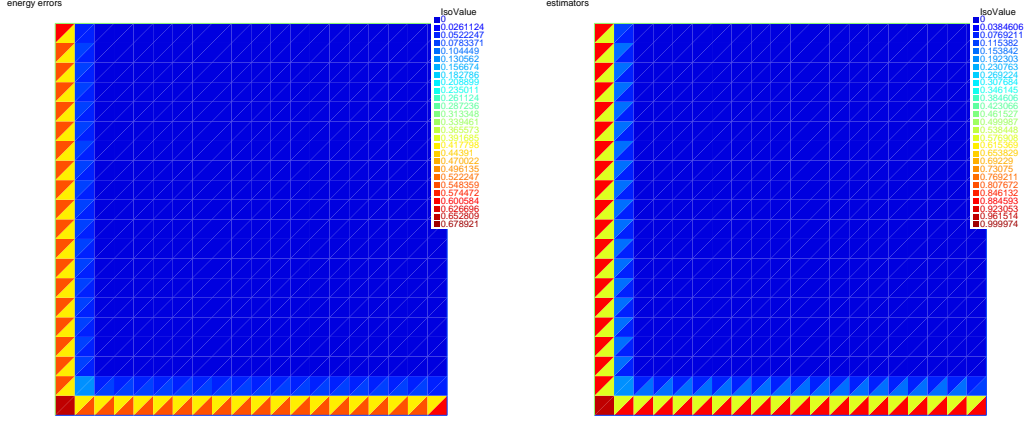


Figure 3: [Section 6.1,  $\varepsilon = 1$ ,  $\kappa = 10^2$ ] Exact (left) and estimated (right) energy error on each mesh element, uniform  $20 \times 20 \times 2$  mesh, polynomial degree  $p = 2$ .

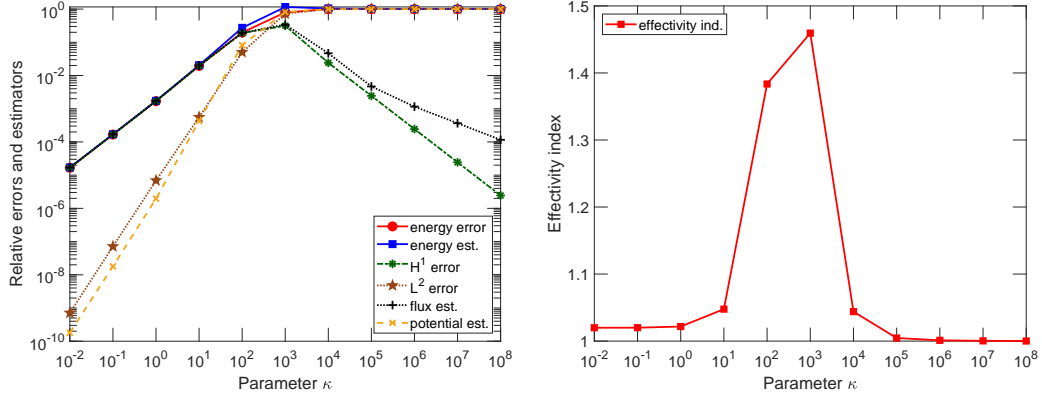


Figure 4: [Section 6.1,  $\varepsilon = 1$ ,  $\kappa$  varies] Energy errors and estimates together with their components, scaled relative to  $\|u_{\mathcal{T}}\|$  (left), and corresponding effectivity indices (right). Uniform  $100 \times 100 \times 2$  mesh, polynomial degree  $p = 1$ .

## 6 Numerical illustration

We illustrate here our theoretical developments on two test cases performed with the FreeFem++ code [24].

### 6.1 A boundary layer

Consider problem (1.1) on the unit square  $\Omega = (0, 1) \times (0, 1)$ , with the exact solution as in [23] given by

$$u(x, y) = e^{-\frac{\kappa}{\varepsilon}x} + e^{-\frac{\kappa}{\varepsilon}y};$$

this corresponds to taking  $f = 0$  in (1.1) but replacing the homogeneous Dirichlet boundary condition by the value of  $u$  on the boundary  $\partial\Omega$ . We take the diffusion parameter  $\varepsilon = 1$  and observe that  $u$  develops a sharp boundary layer along the axes  $y = 0$  and  $x = 0$  for high values of the reaction parameter  $\kappa$ . We employ the a posteriori error estimators of Theorem 3.1, with in

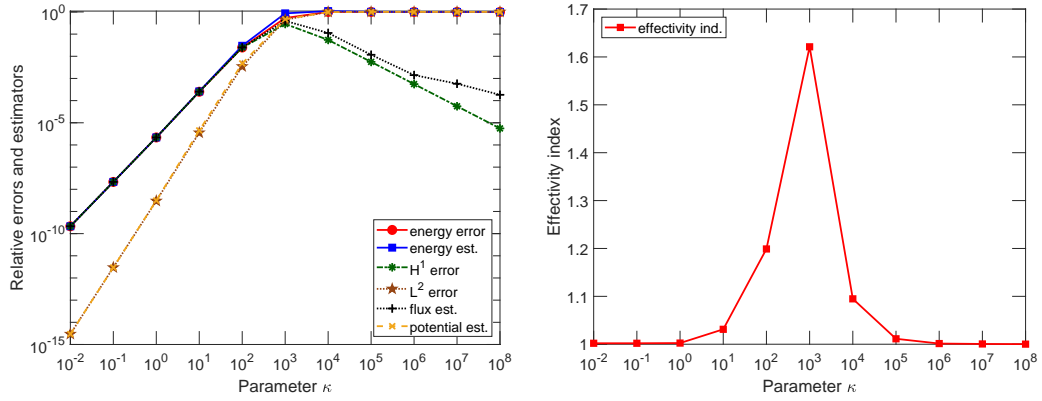


Figure 5: [Section 6.1,  $\varepsilon = 1$ ,  $\kappa$  varies] Energy errors and estimates together with their components, scaled relative to  $\|u_{\mathcal{T}}\|$  (left), and corresponding effectivity indices (right). Uniform  $100 \times 100 \times 2$  mesh, polynomial degree  $p = 2$ .

particular the flux  $\sigma_{\mathcal{T}}$  and the potential  $\phi_{\mathcal{T}}$  constructed following Definition 2.3. The polynomial degree of the finite element approximation (1.5) is equal to either  $p = 1$  or  $p = 2$ ; correspondingly,  $\sigma_{\mathcal{T}}$  is constructed in the  $\mathbf{RTN}_p$  finite-dimensional subspace of  $\mathbf{H}(\text{div}, \Omega)$ , whereas the potential  $\phi_{\mathcal{T}}$  is a piecewise affine or a piecewise quadratic polynomial.

Figure 3 presents the exact element-wise errors  $\|u - u_{\mathcal{T}}\|_K$  together with the element estimators  $\eta_K := w_K \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K + \|\kappa(u_{\mathcal{T}} - \phi_{\mathcal{T}})\|_K$  for  $\kappa = 10^2$  and  $p = 2$  on a uniform mesh given by  $20 \times 20$  squares, each cut into two triangles. We see that the error is concentrated on the boundary layers of the solution, and we observe an excellent match between the exact element-wise errors and the a posteriori estimators.

Figures 4 and 5 then assess the quality of the estimators for  $\kappa$  varying between  $10^{-2}$  and  $10^8$ , uniform  $100 \times 100$  ( $\times 2$ ) meshes, and respectively  $p = 1$  and  $p = 2$ . We observe a stable (and excellent) effectivity index given by the ratio of the estimate of (3.1) over the error  $\|u - u_{\mathcal{T}}\|$ . For a better insight, we compare the  $H^1$ -seminorm part of the error, given by  $\varepsilon \|\nabla(u - u_{\mathcal{T}})\|$ , with the part of the estimator involving the fluxes, i.e.  $(\sum_{K \in \mathcal{T}} [w_K \|\varepsilon \nabla u_{\mathcal{T}} + \varepsilon^{-1} \sigma_{\mathcal{T}}\|_K]^2)^{1/2}$ . We also compare the  $L^2$ -norm part of the error, given by  $\kappa \|u - u_{\mathcal{T}}\|$ , with the part of the estimator involving the potentials, i.e.  $\kappa \|u_{\mathcal{T}} - \phi_{\mathcal{T}}\|$ . We observe that for smaller values of  $\kappa$ , the  $H^1$ -part of the error and flux-part of the estimator dominate, whereas the situation reverses for higher values of  $\kappa$ . Our estimates also appear to predict quite closely these two components of the error.

## 6.2 Necessity of the weights $w_K$

Our second test case corresponds to a two-dimensional analogue of Example 5.2. We consider  $\Omega := (-1/2, 1/2) \times (-1/2, 1/2)$ , along with homogeneous Dirichlet conditions on left and right edges of  $\partial\Omega$ , and homogeneous Neumann boundary conditions on top and bottom edges of  $\partial\Omega$ . We define  $f$  as the extension by constants along the lines  $x = \text{const}$  of the function from (5.5), where we use the value  $m = 3$  for the parameter in (5.5). Then the exact solution  $u$  is simply the extension of the one-dimensional one, see Figure 2. We construct the meshes in analogy with Example 5.2, so the value  $m = 3$  leads to a uniform mesh of  $\Omega$  with  $16 \times 16 \times 2$  elements. Taking the Dirichlet boundary conditions prescribed by the extensions of (5.6) everywhere on  $\partial\Omega$ , the finite element solution (1.5) coincides with (5.6). Figure 6 illustrates this in terms of the absolute values of the pointwise differences  $\varepsilon \partial_x(u - u_{\mathcal{T}})$  and  $\kappa(u - u_{\mathcal{T}})$ .

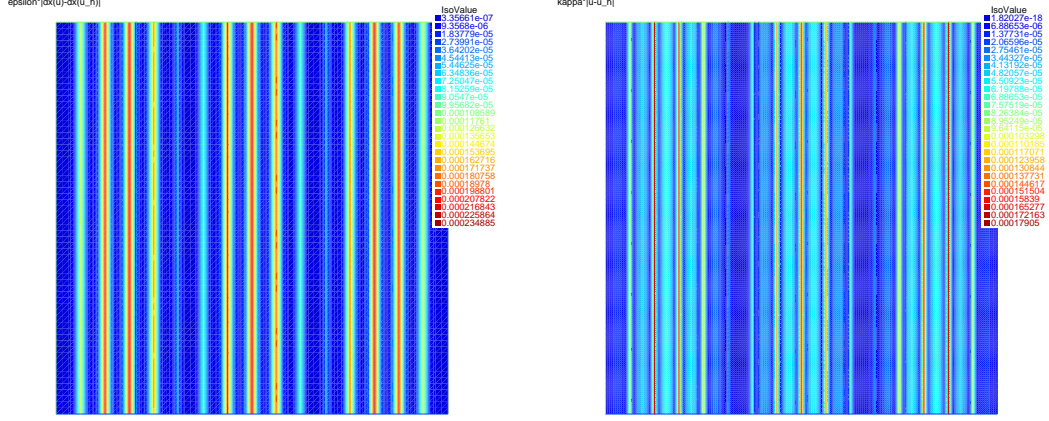


Figure 6: [Section 6.2,  $\varepsilon = 1$ ,  $\kappa = 10^2$ ,  $m = 3$ ] Absolute values of the pointwise differences  $\varepsilon \partial_x(u - u_{\mathcal{T}})$  (left) and  $\kappa(u - u_{\mathcal{T}})$  (right), uniform  $16 \times 16 \times 2$  mesh, polynomial degree  $p = 1$ .

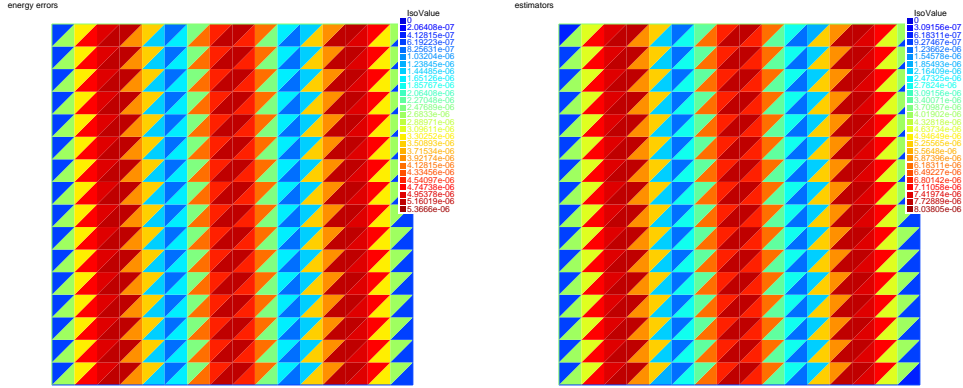


Figure 7: [Section 6.2,  $\varepsilon = 1$ ,  $\kappa = 10^2$ ,  $m = 3$ ] Exact (left) and estimated (right) energy error on each mesh element, uniform  $16 \times 16 \times 2$  mesh, polynomial degree  $p = 1$ .

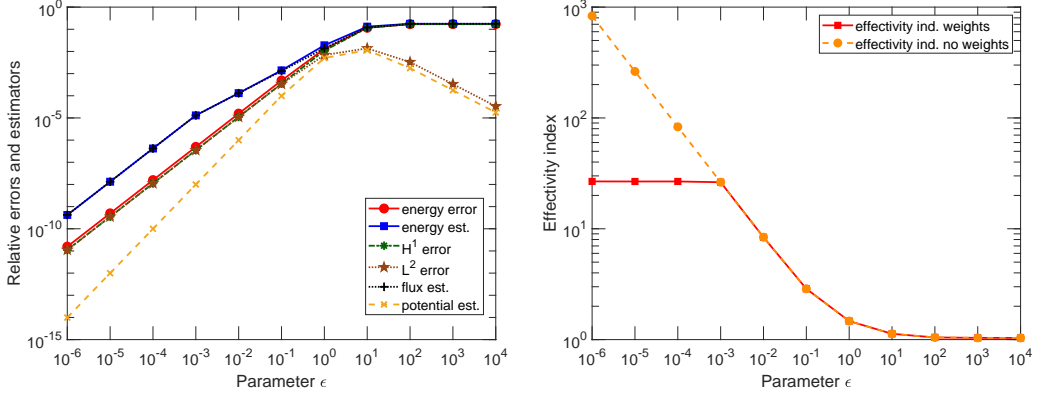


Figure 8: [Section 6.2,  $\varepsilon$  varies,  $\kappa = 10^2$ ,  $m = 3$ ] Relative error and estimate together with their components (left), effectivity indices for weighted and unweighted error estimators (right). Uniform  $16 \times 16 \times 2$  mesh, polynomial degree  $p = 1$ .

We present in Figure 7 the energy errors  $\|u - u_T\|_K$  together with the element estimators  $\eta_K$ , still for  $\varepsilon = 1$ ,  $\kappa = 10^2$ , and  $p = 1$ . The estimators match the true error distribution over  $\Omega$  nearly perfectly, and they identify well the important jumps of the normal component of the approximate solution.

Finally, Figure 8 assesses the quality of our estimates for  $\varepsilon$  varying between  $10^{-6}$  and  $10^4$ ,  $\kappa = 10^2$ , and  $p = 1$ . The effectivity index remains uniformly bounded in all cases, bounded from above by approximately 27, and tends to one for large values of  $\varepsilon$ . The effectivity indices of roughly 27 suggest that the constant  $C_*$  from (2.7) could ideally be reduced, possibly through numerical computation of the optimal constants appearing in the definition of  $C_*$ . The right panel of Figure 8 additionally shows the effectivity indices of the estimators when not employing the weights (neither in Definition 2.3, nor in Theorem 3.1), which become unbounded for small  $\varepsilon$ . Thus, in confirmation of our theory, the weights are crucial for the robustness of the estimators with respect to the ratio  $h\kappa/\varepsilon$ .

## A Explicit constants for the inverse inequality

For each polynomial degree  $p \geq 0$ , let  $C_{p,1}$  denote the best constant of the inverse inequality for the unit interval  $(0, 1)$ , i.e.

$$\|v'\|_{L^2(0,1)} \leq C_{p,1} \|v\|_{L^2(0,1)} \quad \forall v \in \mathbb{P}_p(0, 1), \quad (\text{A.1})$$

where  $\mathbb{P}_p(0, 1)$  denotes the space of univariate polynomials of degree at most  $p$  on  $(0, 1)$ . It was shown in [28] that, for all  $p \geq 0$ ,

$$C_{p,1} \leq \frac{1}{\sqrt{2}} \sqrt{p(p+1)(p+2)(p+3)}; \quad (\text{A.2})$$

here we have taken into account the fact that we consider  $C_{p,1}$  on the unit interval  $(0, 1)$  rather than the interval  $(-1, 1)$  as in [28]. This improves on earlier bounds, e.g. in [32].

We will show here explicit bounds for the constants of the inverse inequality for hypercubes and simplices in terms of  $C_{p,1}$ .

## A.1 Unit hypercube

For an integer  $d \geq 1$ , let  $\{1:d\}$  be a shorthand notation for  $\{1, \dots, d\}$ . Let  $Q_d := \{x \in \mathbb{R}^d, |x|_\infty \leq 1, x_i \geq 0 \ \forall i \in \{1:d\}\}$  denote the unit hypercube in  $\mathbb{R}^d$ , where  $|x|_\infty := \max_{i \in \{1:d\}} |x_i|$ . Let  $\mathbb{P}_p(Q_d)$  denote the space of polynomials of total degree at most  $p$  on  $Q_d$ .

**Lemma A.1.** *For all  $d \geq 1$  and all  $p \geq 0$ , we have*

$$\|v_{x_i}\|_{L^2(Q_d)} \leq C_{p,1} \|v\|_{L^2(Q_d)} \quad \forall v \in \mathbb{P}_p(Q_d), \quad \forall i \in \{1:d\}. \quad (\text{A.3})$$

*Proof.* After a possible re-labelling of the indices, it is enough to show that (A.3) holds for the case  $i = 1$ . Then, writing  $x = (x_1, x')$  with  $x' \in \mathbb{R}^{d-1}$ , we see that

$$\int_{Q^d} |v_{x_1}|^2 dx = \int_{Q^{d-1}} \int_0^1 |v_{x_1}(x_1, x')|^2 dx_1 dx' \leq C_{p,1}^2 \int_{Q^{d-1}} \int_0^1 |v(x_1, x')|^2 dx_1 dx' = C_{p,1}^2 \int_{Q^d} |v|^2 dx,$$

where we use the fact that  $x_1 \mapsto v(x_1, x')$  is in  $\mathbb{P}_p(0, 1)$  for all  $x'$ .  $\square$

## A.2 Unit simplex

For a parameter  $t > 0$ , let  $K_t^d := \{x \in \mathbb{R}^d, |x|_1 \leq t, x_i \geq 0 \ \forall i \in \{1:d\}\}$ , where  $|x|_1 := \sum_{i=1}^d |x_i|$ , denote the simplex in  $\mathbb{R}^d$  with side-length  $t$ . If  $t = 1$ , we adopt the simpler notation  $K^d := K_1^d$ . Let  $C_{p,d}$  denote the best constant such that

$$\|v_{x_i}\|_{L^2(K^d)} \leq C_{p,d} \|v\|_{L^2(K^d)} \quad \forall v \in \mathbb{P}_p(K^d), \quad \forall i \in \{1:d\}. \quad (\text{A.4})$$

We shall obtain here an explicit bound for the constant  $C_{p,d}$  in terms of the space dimension  $d$  and the constant  $C_{p,1}$  of (A.1).

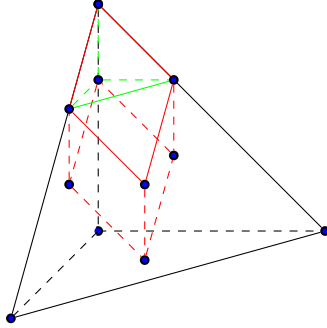


Figure 9: Subdivision of the unit simplex used in the proof of Theorem A.2. The unit simplex is shown for  $d = 3$ , along with its sub-simplex  $K_\dagger$  (edges shown in green) and sub-parallelepipiped  $Q_\dagger$  (edges shown in red).

**Theorem A.2.** *For all  $d \geq 1$  and for all  $p \geq 0$ , the best constant  $C_{p,d}$  in (A.4) satisfies*

$$C_{p,d} \leq \frac{\sqrt{5}}{4} (2\sqrt{2})^d C_{p,1}. \quad (\text{A.5})$$

*Proof.* The proof is based on an induction on the dimension, where we seek to bound  $C_{p,d}$  in terms of  $C_{p,d-1}$ ,  $C_{p,1}$ , and  $d$ . Without loss of generality, it is enough to consider only the case  $i = 1$  in (A.4), after a possible re-labelling of the indices. Then, writing  $x = (x', x_d)$  with  $x' \in \mathbb{R}^{d-1}$ , we have  $\int_{K^d} |v_{x_1}|^2 dx = \int_0^1 \int_{K_{1-x_d}^{d-1}} |v_{x_1}|^2 dx' dx_d$ . Since, for fixed  $x_d \in (0, 1)$ ,  $x' \mapsto v(x', x_d)$  is a polynomial of degree at most  $p$  on  $K_{1-x_d}^{d-1}$ , it would be natural to apply the inverse inequality for simplices of dimension  $d-1$  after a suitable scaling. However, a difficulty arises for  $x_d$  close to 1 due to the appearance of a negative power of  $1 - x_d$  inside the resulting integral. We can overcome this obstacle using an appropriate subdivision of the unit simplex and a change of variables.

The proof proceeds in two steps. We first treat the case  $d = 2$  and show that (A.5) holds (we actually consider  $d \geq 2$  below for the sake of generality), and then the induction is carried out on  $d$  with a different argument, leading to a sharper bound than that would result from step 1 only.

*Step 1.* Let  $d \geq 2$  and consider the partition of  $K$  into  $K_* := \{x \in K, x_d < 1 - 1/d\}$  and  $K_\dagger := K \setminus K_*$ . Then,  $\|v_{x_1}\|_{L^2(K^d)}^2 = \|v_{x_1}\|_{L^2(K_*)}^2 + \|v_{x_1}\|_{L^2(K_\dagger)}^2$ , and the first term can be bounded as follows:

$$\begin{aligned} \|v_{x_1}\|_{L^2(K_*)}^2 &= \int_0^{1-1/d} \left( \int_{K_{1-x_d}^{d-1}} |v_{x_1}|^2 dx' \right) dx_d \\ &\leq \int_0^{1-1/d} \left( \frac{C_{p,d-1}^2}{(1-x_d)^2} \int_{K_{1-x_d}^{d-1}} |v|^2 dx' \right) dx_d \leq d^2 C_{p,d-1}^2 \|v\|_{L^2(K_*)}^2, \end{aligned} \quad (\text{A.6})$$

where crucially we use the fact that  $(1 - x_d)^{-2} \leq d^2$  for  $x_d \leq 1 - 1/d$ . In order to bound the second term  $\|v_{x_1}\|_{L^2(K_\dagger)}^2$ , we introduce a change of coordinates in terms of the affine map  $F$  defined by

$$F(\xi) := e_d + \sum_{i=1}^d (e_{i-1} - e_d) \xi_i,$$

where  $e_0 = 0$ , and  $e_i$  is the  $i$ -th unit vector for  $1 \leq i \leq d$ . Letting  $x = F(\xi)$ , we have  $x_j = \xi_{j+1}$  for  $j \leq d-1$ , and  $x_d = 1 - \sum_{i=1}^d \xi_i$ . The inverse is then given by  $\xi_1 = 1 - \sum_{i=1}^d x_i$ , and  $\xi_j = x_{j-1}$  for  $2 \leq j \leq d$ . It is thus easily seen that  $F$  is a bijection from  $K$  onto itself, and that  $F(0) = x_d$ . Thus  $F$  corresponds to a change of coordinates on the unit simplex. Additionally, it can be shown that the Jacobian  $|\det DF| = 1$ .

Let  $Q_{1/d}^d := \{\xi \in Q^d, |\xi|_\infty \leq 1/d\}$  be a hypercube with side length  $1/d$ , and let  $Q_\dagger$  be the parallelepiped obtained as the image of  $Q_{1/d}^d$  under the mapping  $F$ , i.e.  $Q_\dagger = F(Q_{1/d}^d)$ . It is then easy, but tedious, to show that

$$K_\dagger \subset Q_\dagger \subset K. \quad (\text{A.7})$$

Figure 9 illustrates the sets  $K_\dagger$ ,  $Q_\dagger$ , and  $K$  for the case  $d = 3$ . Now, let  $\tilde{v}(\xi) = v(F(\xi))$  be the pullback of  $v$  under  $F$ . Since  $F$  is affine,  $\tilde{v} \in \mathbb{P}_p(K^d)$ . It is also easy to check that  $v_{x_1} = \tilde{v}_{\xi_2} - \tilde{v}_{\xi_1}$ . Using the change of variables and the fact that  $|\det DF| = 1$ , it follows from (A.7) that

$$\|v_{x_1}\|_{L^2(K_\dagger)}^2 \leq \|v_{x_1}\|_{L^2(Q_\dagger)}^2 = \|\tilde{v}_{\xi_2} - \tilde{v}_{\xi_1}\|_{L^2(Q_{1/d}^d)}^2 \leq 2(\|\tilde{v}_{\xi_2}\|_{L^2(Q_{1/d}^d)}^2 + \|\tilde{v}_{\xi_1}\|_{L^2(Q_{1/d}^d)}^2).$$

Applying the inverse inequality for hypercubes, namely  $\|\tilde{v}_{\xi_i}\|_{L^2(Q_{1/d}^d)}^2 \leq d^2 C_{p,1}^2 \|\tilde{v}\|_{L^2(Q_{1/d}^d)}^2$ , and changing back to the original variables, we then obtain from the second inclusion in (A.7) that

$$\|v_{x_1}\|_{L^2(K_\dagger)}^2 \leq 4d^2 C_{p,1}^2 \|v\|_{L^2(Q_\dagger)}^2 \leq 4d^2 C_{p,1}^2 \|v\|_{L^2(K^d)}^2. \quad (\text{A.8})$$

Therefore, combining (A.6) and (A.8), we arrive at  $\|v_{x_1}\|_{L^2(K^d)}^2 \leq d^2(C_{p,d-1}^2 + 4C_{p,1}^2)\|v\|_{L^2(K^d)}^2$ , for any  $v \in \mathbb{P}_p(K^d)$ . This implies  $C_{p,d}^2 \leq d^2(C_{p,d-1}^2 + 4C_{p,1}^2)$ , and thus (A.1) and an induction argument show that

$$C_{p,d} \leq \left(1 + 4 \sum_{j=1}^{d-1} \frac{1}{(j!)^2}\right)^{\frac{1}{2}} d! C_{p,1} \quad (\text{A.9})$$

for any  $d \geq 2$ . This shows (A.4), but with a worse constant than that of (A.5) for  $d \geq 3$ . For this reason, we proceed in a second step in a different way.

*Step 2.* Let  $d \geq 3$ . We again subdivide the simplex  $K$ , this time as

$$K = \{x \in K, x_d \leq 1/2\} \cup \{x \in K, x_{d-1} \leq 1/2\}.$$

Furthermore, for any fixed  $x_{d-1}, x'_{d-1} = (x_1, \dots, x_{d-2}, x_d)' \mapsto v(x)$  is a polynomial of degree at most  $p$  on a simplex that is isometric to  $K_{1-x_{d-1}}^{d-1}$ . Let also  $x'_d = (x_1, \dots, x_{d-1})$ . Crucially, since  $d \geq 3$  and we subdivide above into two subsets, we can avoid the critical subset  $K_{\dagger}$  of Step 1 as

$$\begin{aligned} \|v_{x_1}\|_{L^2(K^d)}^2 &\leq \sum_{j=d-1}^d \int_0^{1/2} \left( \int_{K_{1-x_j}^{d-1}} |v_{x_1}|^2 dx'_j \right) dx_j \\ &\leq \sum_{j=d-1}^d \int_0^{1/2} \frac{C_{p,d-1}^2}{(1-x_j)^2} \left( \int_{K_{1-x_j}^{d-1}} |v|^2 dx'_j \right) dx_j \leq 8C_{p,d-1}^2 \|v\|_{K^d}^2. \end{aligned} \quad (\text{A.10})$$

It then follows by induction that  $C_{p,d} \leq (2\sqrt{2})^{d-2} C_{p,2}$  for all  $d \geq 3$ . Since  $C_{p,2} \leq 2\sqrt{5} C_{p,1}$  by (A.9), we get (A.5).  $\square$

Applying Theorem A.2 to the cases  $d = 2$  and  $d = 3$  gives the following explicit bounds

$$C_{p,2} \leq \sqrt{10p(p+1)(p+2)(p+3)}, \quad C_{p,3} \leq \sqrt{80p(p+1)(p+2)(p+3)}. \quad (\text{A.11})$$

### A.3 General simplex

Let  $K$  be a simplex in  $\mathbb{R}^d$ ,  $d \geq 2$ , and let  $\hat{K}$  denote the unit simplex. Let  $J_K$  denote the differential of the affine transformation mapping  $T_K: \hat{K} \rightarrow K$ . For  $\mathbf{v} \in \mathbf{RTN}_p(K)$ , we define the Piola transformation  $\hat{\mathbf{v}} \in \mathbf{RTN}_p(\hat{K})$  by

$$\hat{\mathbf{v}}(\hat{x}) = |\det J_K| J_K^{-1} [\mathbf{v} \circ T_K(\hat{x})]. \quad (\text{A.12})$$

**Lemma A.3** (Ciarlet [11] Thm 3.1.2 and [14]). *There holds*

$$\|J_K\|_2 \leq \frac{h_K}{\rho_{\hat{K}}}, \quad \|J_K^{-1}\|_2 \leq \frac{\sqrt{2}}{\rho_{\hat{K}}}, \quad |\det J_K| = \frac{|K|_d}{|\hat{K}|_d}, \quad \frac{|\partial K|_{d-1}}{|K|_d} \leq (d+1)d\vartheta_{\mathcal{T}} h_K^{-1}. \quad (\text{A.13})$$

Note that in Lemma A.3, we have used the fact that the diameter of the unit simplex is  $\sqrt{2}$  for all  $d \geq 2$ .

**Lemma A.4** (Warburton & Hesthaven [40]). *Let  $v \in \mathbb{P}_p(K)$ . Then*

$$\|v\|_{\partial K} \leq \sqrt{\frac{(p+1)(p+d)}{d} \frac{|\partial K|_{d-1}}{|K|_d}} \|v\|_K.$$

Therefore, for  $\mathbf{v} \in \mathbf{RTN}_p(K)$ , we have

$$h_K^{1/2} \|\mathbf{v} \cdot \mathbf{n}\|_{\partial K} \leq C_{\text{inv},p,\partial} \|\mathbf{v}\|_K, \quad C_{\text{inv},p,\partial} := \sqrt{(d+1)(p+2)(p+d+1)\vartheta_{\mathcal{T}}}.$$

**Lemma A.5.** *Let  $K$  be a simplex in  $\mathbb{R}^d$  and  $\mathbf{v} \in \mathbf{RTN}_p(K)$ . Then,*

$$h_K \|\nabla \cdot \mathbf{v}\|_K \leq C_{\text{inv},p} \|\mathbf{v}\|_K, \quad C_{\text{inv},p} := \sqrt{2d} \vartheta_{\mathcal{T}} C_{p+1,d}, \quad (\text{A.14})$$

where  $C_{p,d}$  is characterized in Theorem A.2.

*Proof.* Using the Piola Transformation, we have  $\nabla \cdot \mathbf{v} = \nabla_{\hat{x}} \cdot \hat{\mathbf{v}} / |\det J_K|$ , therefore,  $\|\nabla \cdot \mathbf{v}\|_K^2 \leq |\det J_K|^{-1} \|\nabla \cdot \hat{\mathbf{v}}\|_{\hat{K}}^2$ . Then, since  $\hat{\mathbf{v}}_i \in \mathbb{P}_{p+1}(\hat{K})$  for  $i \in \{1:d\}$ , we apply Theorem A.2 to obtain  $\|\nabla \cdot \hat{\mathbf{v}}\|_{\hat{K}}^2 \leq d \sum_{i=1}^d \|\hat{\mathbf{v}}_{i,x_i}\|_{\hat{K}}^2 \leq d C_{p+1,d}^2 \|\hat{\mathbf{v}}\|_{\hat{K}}^2$ . Then, using the definition of the Piola transformation, it is seen that  $\|\hat{\mathbf{v}}\|_{\hat{K}}^2 \leq \|J_K^{-1}\|_2^2 |\det J_K| \|\mathbf{v}\|_K^2$ . We then use the bound  $\|J_K^{-1}\|_2 \leq \sqrt{2} \vartheta_{\mathcal{T}} h_K^{-1}$  from (A.13) to find that  $h_K \|\nabla \cdot \mathbf{v}\|_K \leq \sqrt{2d} \vartheta_{\mathcal{T}} C_{p+1,d} \|\mathbf{v}\|_K$ , which finishes the proof.  $\square$

## References

- [1] M. AINSWORTH, A. ALLENDES, G. R. BARRENECHEA, AND R. RANKIN, *Fully computable a posteriori error bounds for stabilised FEM approximations of convection–reaction–diffusion problems in three dimensions*, Internat. J. Numer. Methods Fluids, 73 (2013), pp. 765–790.
- [2] M. AINSWORTH AND I. BABUŠKA, *Reliable and robust a posteriori error estimation for singularly perturbed reaction-diffusion problems*, SIAM J. Numer. Anal., 36 (1999), pp. 331–353.
- [3] M. AINSWORTH AND J. T. ODEN, *A posteriori error estimation in finite element analysis*, Pure and Applied Mathematics (New York), Wiley-Interscience [John Wiley & Sons], New York, 2000.
- [4] M. AINSWORTH AND T. VEJCHODSKÝ, *Fully computable robust a posteriori error bounds for singularly perturbed reaction-diffusion problems*, Numer. Math., 119 (2011), pp. 219–243.
- [5] ———, *Robust error bounds for finite element approximation of reaction-diffusion problems with non-constant reaction coefficient in arbitrary space dimension*, Comput. Methods Appl. Mech. Engrg., 281 (2014), pp. 184–199.
- [6] ———, *A simple approach to reliable and robust a posteriori error estimation for singularly perturbed problems*, Comput. Methods Appl. Mech. Engrg., 353 (2019), pp. 373–390.
- [7] T. APEL, S. NICAISE, AND D. SIRCH, *A posteriori error estimation of residual type for anisotropic diffusion-convection-reaction problems*, J. Comput. Appl. Math., 235 (2011), pp. 2805–2820.
- [8] D. BRAESS, V. PILLWEIN, AND J. SCHÖBERL, *Equilibrated residual error estimates are p-robust*, Comput. Methods Appl. Mech. Engrg., 198 (2009), pp. 1189–1197.
- [9] D. BRAESS AND J. SCHÖBERL, *Equilibrated residual error estimator for edge elements*, Math. Comp., 77 (2008), pp. 651–672.
- [10] I. CHEDDADI, R. FUČÍK, M. I. PRIETO, AND M. VOHRALÍK, *Guaranteed and robust a posteriori error estimates for singularly perturbed reaction-diffusion problems*, M2AN Math. Model. Numer. Anal., 43 (2009), pp. 867–888.



- [11] P. G. CIARLET, *The Finite Element Method for Elliptic Problems*, vol. 4 of Studies in Mathematics and its Applications, North-Holland, Amsterdam, 1978.
- [12] A. DEMLOW AND N. KOPTEVA, *Maximum-norm a posteriori error estimates for singularly perturbed elliptic reaction-diffusion problems*, Numer. Math., 133 (2016), pp. 707–742.
- [13] P. DESTUYNDER AND B. MÉTIVET, *Explicit error bounds in a conforming finite element method*, Math. Comp., 68 (1999), pp. 1379–1396.
- [14] D. A. DI PIETRO AND A. ERN, *Mathematical aspects of discontinuous Galerkin methods*, vol. 69 of Mathématiques & Applications (Berlin) [Mathematics & Applications], Springer, Heidelberg, 2012.
- [15] V. DOLEJŠÍ, A. ERN, AND M. VOHRALÍK, *hp-adaptation driven by polynomial-degree-robust a posteriori error estimates for elliptic problems*, SIAM J. Sci. Comput., 38 (2016), pp. A3220–A3246.
- [16] M. EIGEL AND C. MERDON, *Equilibration a posteriori error estimation for convection-diffusion-reaction problems*, J. Sci. Comput., 67 (2016), pp. 747–768.
- [17] M. EIGEL AND T. SAMROWSKI, *Functional a posteriori error estimation for stationary reaction-convection-diffusion problems*, Comput. Methods Appl. Math., 14 (2014), pp. 135–150.
- [18] A. ERN, I. SMEARS, AND M. VOHRALÍK, *Discrete  $p$ -robust  $\mathbf{H}(\text{div})$ -liftings and a posteriori estimates for elliptic problems with  $H^{-1}$  source terms*, Calcolo, 54 (2017), pp. 1009–1025.
- [19] A. ERN AND M. VOHRALÍK, *Adaptive inexact Newton methods with a posteriori stopping criteria for nonlinear diffusion PDEs*, SIAM J. Sci. Comput., 35 (2013), pp. A1761–A1791.
- [20] ———, *Polynomial-degree-robust a posteriori estimates in a unified setting for conforming, nonconforming, discontinuous Galerkin, and mixed discretizations*, SIAM J. Numer. Anal., 53 (2015), pp. 1058–1081.
- [21] ———, *Stable broken  $H^1$  and  $\mathbf{H}(\text{div})$  polynomial extensions for polynomial-degree-robust potential and flux reconstruction in three space dimensions*, Math. Comp., 89 (2020), pp. 551–594.
- [22] M. FAUSTMANN AND J. M. MELENK, *Robust exponential convergence of hp-FEM in balanced norms for singularly perturbed reaction-diffusion problems: corner domains*, Comput. Math. Appl., 74 (2017), pp. 1576–1589.
- [23] S. GROSMAN, *An equilibrated residual method with a computable error approximation for a singularly perturbed reaction-diffusion problem on anisotropic finite element meshes*, M2AN Math. Model. Numer. Anal., 40 (2006), pp. 239–267.
- [24] F. HECHT, *New development in FreeFem++*, J. Numer. Math., 20 (2012), pp. 251–265.
- [25] N. KOPTEVA, *Maximum-norm a posteriori error estimates for singularly perturbed reaction-diffusion problems on anisotropic meshes*, SIAM J. Numer. Anal., 53 (2015), pp. 2519–2544.
- [26] ———, *Energy-norm a posteriori error estimates for singularly perturbed reaction-diffusion problems on anisotropic meshes*, Numer. Math., 137 (2017), pp. 607–642.

- [27] ———, *Fully computable a posteriori error estimator using anisotropic flux equilibration on anisotropic meshes*. ArXiv Preprint 1704.04404, 2017.
- [28] C. KOUTSCHAN, M. NEUMÜLLER, AND C.-S. RADU, *Inverse inequality estimates with symbolic computation*, Adv. in Appl. Math., 80 (2016), pp. 1–23.
- [29] G. KUNERT, *Robust a posteriori error estimation for a singularly perturbed reaction-diffusion equation on anisotropic tetrahedral meshes*, Adv. Comput. Math., 15 (2001), pp. 237–259. A posteriori error estimation and adaptive computational methods.
- [30] T. LINSS, *A posteriori error estimation for arbitrary order FEM applied to singularly perturbed one-dimensional reaction-diffusion problems*, Appl. Math., 59 (2014), pp. 241–256.
- [31] S. I. REPIN AND S. SAUTER, *Functional a posteriori estimates for the reaction-diffusion problem*, C. R. Math. Acad. Sci. Paris, 343 (2006), pp. 349–354.
- [32] C. SCHWAB, *p- and hp-finite element methods*, Numerical Mathematics and Scientific Computation, The Clarendon Press, Oxford University Press, New York, 1998. Theory and applications in solid and fluid mechanics.
- [33] R. P. STEVENSON, *The uniform saturation property for a singularly perturbed reaction-diffusion equation*, Numer. Math., 101 (2005), pp. 355–379.
- [34] T. VEJCHODSKÝ, *Complementarity based a posteriori error estimates and their properties*, Math. Comput. Simulation, 82 (2012), pp. 2033–2046.
- [35] ———, *On the quality of local flux reconstructions for guaranteed error bounds*, in Applications of mathematics 2015, Czech. Acad. Sci., Prague, 2015, pp. 242–255.
- [36] R. VERFÜRTH, *Robust a posteriori error estimators for a singularly perturbed reaction-diffusion equation*, Numer. Math., 78 (1998), pp. 479–493.
- [37] ———, *Robust a posteriori error estimates for stationary convection-diffusion equations*, SIAM J. Numer. Anal., 43 (2005), pp. 1766–1782.
- [38] ———, *A note on constant-free a posteriori error estimates*, SIAM J. Numer. Anal., 47 (2009), pp. 3180–3194.
- [39] ———, *A posteriori error estimation techniques for finite element methods*, Numerical Mathematics and Scientific Computation, Oxford University Press, Oxford, 2013.
- [40] T. WARBURTON AND J. S. HESTHAVEN, *On the constants in hp-finite element trace inverse inequalities*, Comput. Methods Appl. Mech. Engrg., 192 (2003), pp. 2765–2773.
- [41] J. ZHAO AND S. CHEN, *Robust a posteriori error estimates for conforming discretizations of a singularly perturbed reaction-diffusion problem on anisotropic meshes*, Adv. Comput. Math., 40 (2014), pp. 797–818.